
Accelerating Diffusion Models with Parallel Sampling: Inference at Sub-Linear Time Complexity

Haoxuan Chen*

ICME

Stanford University

haoxuanc@stanford.edu

Yinuo Ren[†]

ICME

Stanford University

yinuoren@stanford.edu

Lexing Ying

Department of Mathematics and ICME

Stanford University

lexing@stanford.edu

Grant M. Rotskoff

Department of Chemistry and ICME

Stanford University

rotskoff@stanford.edu

Abstract

Diffusion models have become a leading method for generative modeling of both image and scientific data. As these models are costly to train and *evaluate*, reducing the inference cost for diffusion models remains a major goal. Inspired by the recent empirical success in accelerating diffusion models via the parallel sampling technique [1], we propose to divide the sampling process into $\mathcal{O}(1)$ blocks with parallelizable Picard iterations within each block. Rigorous theoretical analysis reveals that our algorithm achieves $\tilde{\mathcal{O}}(\text{poly } \log d)$ overall time complexity, marking *the first implementation with provable sub-linear complexity w.r.t. the data dimension d* . Our analysis is based on a generalized version of Girsanov’s theorem and is compatible with both the SDE and probability flow ODE implementations. Our results shed light on the potential of fast and efficient sampling of high-dimensional data on fast-evolving modern large-memory GPU clusters.

1 Introduction

Diffusion and probability flow based models [2–11] are now state-of-the-art in many fields, such as computer vision and image generation [12–22], natural language processing [23, 24], audio and video generation [25–29], optimization [30, 31], sampling and learning of fixed classes of distributions [32–41], solving high-dimensional partial differential equations [42–46], and more recently several applications in physical, chemical and biological fields [47–63]. For a more comprehensive list of related work, one may refer to the following review papers [64–66]. While there are already many variants, such as denoising diffusion probabilistic models (DDPMs) [7], score-based generative models (SGMs) [9], diffusion schrödinger bridges [67], stochastic interpolants and flow matching [2–4], *etc.*, the recurring idea is to design a stochastic process that interpolates between the data distribution and some simple distribution, along which *score functions* or alike are learned by neural network-based estimators, and then perform inference guided by the learned score functions.

Due to the sequential nature of the sampling process, the inference of high-quality samples from diffusion models often requires a large number of iterations and, thus, evaluations of the neural network-based score function, which can be computationally expensive [68]. Efforts have been

*Equal contribution, alphabetical order.

[†]Corresponding author.

Work	Implementation	Measure	Approx. Time Complexity
[100, Theorem 2]	SDE	$\text{TV}(p_0, \hat{q}_T)^2$	$\tilde{\mathcal{O}}(d\delta^{-1})$
[104, Theorem 2]	SDE	$D_{\text{KL}}(p_\eta \ \hat{q}_{T-\eta})$	$\tilde{\mathcal{O}}(d^2\delta^{-2})$
[107, Corollary 1]	SDE	$D_{\text{KL}}(p_\eta \ \hat{q}_{T-\eta})$	$\tilde{\mathcal{O}}(d\delta^{-2})$
[111, Theorem 3]	ODE w/UMLC correction	$\text{TV}(p_\eta, \hat{q}_{T-\eta})^2$	$\tilde{\mathcal{O}}(\sqrt{d}\delta^{-1})$
Theorem 3.3	SDE w/parallel sampling	$D_{\text{KL}}(p_\eta \ \hat{q}_{T-\eta})$	$\tilde{\mathcal{O}}(\text{poly log}(d\delta^{-2}))$
Theorem 3.5	ODE w/parallel sampling	$\text{TV}(p_\eta, \hat{q}_{T-\eta})^2$	$\tilde{\mathcal{O}}(\text{poly log}(d\delta^{-2}))$

Table 1: Comparison of the approximate time complexity (*cf.* Definition 2.1) of different implementations of diffusion models. η is a small parameter that controls the smooth approximation of the data distribution (*cf.* Section 3.1.1).

made to accelerate this process by resorting to higher-order or randomized numerical schemes [69–79], augmented dynamics [80], adaptive step sizes [81], operator learning [82], restart sampling [83], self-consistency [84–87] and knowledge distillation [88–90]. Recently, several empirical works [1, 91–94] leverage the Picard iteration and triangular Anderson acceleration to parallelize the sampling procedure of diffusion models and achieve empirical success in large-scale image generation tasks. Some other recent work [95, 96] also combine the parallel sampling technique with the randomized midpoint method [97] to accelerate the inference of diffusion models.

This efficiency issue is closely related to the problem of bounding the required number of steps and evaluations of score functions to approximate an arbitrary data distribution on \mathbb{R}^d to δ -accuracy, which has been analyzed extensively in the literature [98–115]. In terms of the dependency on the dimension d , the current state-of-the-art result for the SDE implementation of diffusion models is $\tilde{\mathcal{O}}(d)$ [107], improved from the previous $\tilde{\mathcal{O}}(d^2)$ bound [104]. [111] gives a $\tilde{\mathcal{O}}(\sqrt{d})$ bound for the probability flow ODE implementation by considering a predictor-corrector scheme with the underdamped Langevin Monte Carlo (UMLC) algorithm.

In this work, we aim to provide parallelization strategies, rigorous analysis, and theoretical guarantees for accelerating the inference process of diffusion models. The time complexity of previous implementations of diffusion models has been largely hindered by the discretization error, which requires the step size to scale with $\tilde{\mathcal{O}}(1/d)$ for the SDE implementation and $\tilde{\mathcal{O}}(1/\sqrt{d})$ for the probability flow ODE implementation. We show that the inference process can be first divided into $\mathcal{O}(1)$ blocks with parallelizable evaluations of the score function within each, and thus reduce the overall time complexity to $\tilde{\mathcal{O}}(\text{poly log } d)$. We provide **the first implementation of diffusion models with poly-logarithmic complexity**, a significant improvement over the current state-of-the-art polynomial results that sheds light on the potential fast and efficient sampling of high-dimensional distributions with diffusion models on fast-developing memory-efficient modern GPU clusters.

1.1 Contributions

- We propose parallelized inference algorithms for diffusion models in both the SDE and probability flow ODE implementations (PIADM-SDE/ODE) with exponential integrators, a shrinking step size scheme towards the data end, and the early stopping technique;
- We provide a rigorous convergence analysis of PIADM-SDE, showing that our parallelization strategy yields a diffusion model with $\tilde{\mathcal{O}}(\text{poly log } d)$ approximate time complexity;
- We show that our strategy is also compatible with the probability flow ODE implementation, and PIADM-ODE could improve the space complexity from $\tilde{\mathcal{O}}(d^2)$ to $\tilde{\mathcal{O}}(d^{3/2})$ while maintaining the poly-logarithmic time complexity.

2 Preliminaries

In this section, we briefly recapitulate the framework of score-based diffusion models, define notations, and discuss related work.

2.1 Diffusion Models

In score-based diffusion models, one considers a diffusion process $(\mathbf{x}_s)_{s \geq 0}$ in \mathbb{R}^d governed by the following stochastic differential equation (SDE):

$$d\mathbf{x}_s = \beta_s(\mathbf{x}_s)ds + \sigma_s d\mathbf{w}_s, \quad \text{with } \mathbf{x}_0 \sim p_0, \quad (2.1)$$

where $(\mathbf{w}_s)_{s \geq 0}$ is a standard Brownian motion, and p_0 is the target distribution that we would like to sample from. The distribution of \mathbf{x}_s is denoted by p_s . Once the drift $\beta_s(\cdot)$, the diffusion coefficient σ_s , and a sufficiently large time horizon T are specified, (2.1) also corresponds to a backward process $(\tilde{\mathbf{x}}_t)_{0 \leq t \leq T}$ for another arbitrary diffusion coefficient $(\mathbf{v}_s)_{s \geq 0}$ [116]:

$$d\tilde{\mathbf{x}}_t = \left[-\tilde{\beta}_t(\tilde{\mathbf{x}}_t) + \frac{\tilde{\sigma}_t \tilde{\sigma}_t^\top + \tilde{\mathbf{v}}_t \tilde{\mathbf{v}}_t^\top}{2} \nabla \log \tilde{p}_t(\tilde{\mathbf{x}}_t) \right] dt + \tilde{\mathbf{v}}_t d\mathbf{w}_t, \quad (2.2)$$

where $\tilde{*}_t$ denotes $*_{T-t}$, with $\tilde{p}_0 = p_T$ and $\tilde{p}_T = p_0$.

For notational simplicity, we adopt a simple choice of the drift and the diffusion coefficients in what follows: $\beta_t(\mathbf{x}) = -\frac{1}{2}\mathbf{x}$, $\sigma_t = \mathbf{I}_d$, and $\mathbf{v} = v\mathbf{I}_d$, under which (2.1) is an Ornstein-Uhlenbeck (OU) process converging exponentially to its stationary distribution, *i.e.* $p_T \approx \hat{p}_T := \mathcal{N}(0, \mathbf{I}_d)$, and (2.1) and (2.2) reduce to the following form:

$$d\mathbf{x}_s = -\frac{1}{2}\mathbf{x}_s ds + d\mathbf{w}_s, \quad \text{and} \quad d\tilde{\mathbf{x}}_t = \left[\frac{1}{2}\tilde{\mathbf{x}}_t + \frac{1+v^2}{2} \nabla \log \tilde{p}_t(\tilde{\mathbf{x}}_t) \right] dt + v d\mathbf{w}_t. \quad (2.3)$$

In practice, the score function $\nabla \tilde{p}_t(\tilde{\mathbf{x}}_t)$ is often estimated by a neural network (NN) $\mathbf{s}_t^\theta(\mathbf{x}_t)$, where θ represents its parameters, by minimizing the denoising score-matching loss [117, 118]:

$$\begin{aligned} \mathcal{L}(\theta) &:= \mathbb{E}_{\mathbf{x}_t \sim p_t} \left[\left\| \nabla \log p_t(\mathbf{x}_t) - \mathbf{s}_t^\theta(\mathbf{x}_t) \right\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim p_0} \left[\mathbb{E}_{\mathbf{x}_t \sim p_t | \mathbf{x}_0} \left[\left\| \frac{\mathbf{x}_t - \mathbf{x}_0 e^{-t/2}}{1 - e^{-t}} - \mathbf{s}_t^\theta(\mathbf{x}_t) \right\|^2 \right] \right], \end{aligned} \quad (2.4)$$

and the backward process in (2.3) is approximated by the following SDE thereafter:

$$d\mathbf{y}_t = \left[\frac{1}{2}\mathbf{y}_t + \frac{1+v^2}{2} \mathbf{s}_t^\theta(\mathbf{y}_t) \right] dt + v d\mathbf{w}_t, \quad \text{with } \mathbf{y}_0 \sim \mathcal{N}(0, \mathbf{I}_d). \quad (2.5)$$

Implementations. Diffusion models admit multiple *implementations* depending on the choice of the parameter v in the backward process (2.2). The SDE implementation with $v = 1$ is widely used in the literature for its simplicity and efficiency [10], while recent studies [111] claim that the probability flow ODE implementation with $v = 0$ may exhibit better time complexity. We refer to [111, 119] for theoretical and [120, 121] for empirical comparisons of different implementations.

2.2 Parallel Sampling

Parallel sampling algorithms have been actively explored in the literature, including the parallel tempering method [122–124] and several recent studies [125–127]. For diffusion models, the idea of parallel sampling is based on the *Picard iteration* [128, 129] for solving nonlinear ODEs. Suppose we have an ODE $d\mathbf{x}_t = \mathbf{f}_t(\mathbf{x}_t)dt$ and we would like to solve it for $t \in [0, T]$, then the Picard iteration is defined as follows:

$$\mathbf{x}_t^{(0)} \equiv \mathbf{x}_0, \quad \text{and} \quad \mathbf{x}_t^{(k+1)} := \mathbf{x}_0 + \int_0^t \mathbf{f}_s(\mathbf{x}_s^{(k)})ds, \quad \text{for } k \in [0 : K-1]. \quad (2.6)$$

Under assumptions on the Lipschitz continuity of \mathbf{f}_t , the Picard iteration converges to the true solution exponentially fast, in the sense that $\|\mathbf{x}_t^{(k)} - \mathbf{x}_t\|_{L^\infty([0, T])} \leq \delta$ with $K = \mathcal{O}(\log \delta^{-1})$ iterations. Unlike high-order ODE solvers, the Picard iteration is intrinsically parallelizable: for any $t \in [0, T]$, the computation of $\mathbf{x}_t^{(k+1)}$ relies merely on the values of the most recent iteration $\mathbf{x}_t^{(k)}$. With sufficient computational sources parallelizing the evaluations of \mathbf{f} , the computational cost of solving the ODE no longer scales with T but with the number of iterations K .

Recently, this idea has been applied to both the Langevin Monte Carlo (LMC) and the underdamped Langevin Monte Carlo (ULMC) contexts [130]. Roughly speaking, it is proposed to simulate the Langevin diffusion process $d\mathbf{x}_t = -\nabla V(\mathbf{x}_t)dt + d\mathbf{w}_t$ with the following iteration resembling (2.6):

$$\mathbf{x}_t^{(0)} \equiv \mathbf{x}_0, \quad \text{and} \quad \mathbf{x}_t^{(k+1)} := \mathbf{x}_0 - \int_0^t \nabla V(\mathbf{x}_s^{(k)})ds + \mathbf{w}_t, \quad \text{for } k \in [0 : K - 1], \quad (2.7)$$

where all iterations share a common Wiener process $(\mathbf{w}_t)_{t \geq 0}$.

It is shown that for well-conditioned log-concave distributions, parallelized LMC would achieve an iteration depth of $K = \tilde{O}(\text{poly log } d)$ that matches the indispensable time horizon $T = \tilde{O}(\text{poly log } d)$ to achieve exponential ergodicity (cf. [130, Theorem 13]). This promises a significant speedup in sampling high-dimensional distributions from the standard LMC of $T = \tilde{O}(d)$ iterations, hindered by the $o(1/d)$ step size as imposed by the discretization error and now evaded by the parallelization.

2.3 Approximate Time Complexity

A similar situation is expected in diffusion models, where the application bottleneck is largely the inference process with sequential iterations and expensive evaluations of the learned score function $\mathbf{s}_t^\theta(\cdot)$, which is often parametrized by large-scale NNs. Despite several unavoidable costs involving pre- and post-processing, data storage and retrieval, and arithmetic operations, we define the following notion of the *approximate time complexity* of the inference process of diffusion models:

Definition 2.1 (Approximate time complexity). *For a specific implementation of diffusion models (2.5), we define the approximate time complexity of the sampling process as the number of unparallelizable evaluations of the learned NN-based score function $\mathbf{s}_t^\theta(\cdot)$.*

This definition coincides with the notion of *the number of steps required to reach a certain accuracy* in [104, 100], *iteration complexity* in [107, 111], etc. in the previous theoretical studies of diffusion models. We have adopted this notion in Table 1 for a comparison of the current state-of-the-art results and our bounds in this work. We will use the notion of *space complexity* likewise to denote the approximate required storage during the inference. Trivially, the space complexity of the sequential implementation is $\mathcal{O}(d)$. Should no confusion occur, we omit the dependency of the complexities above on the accuracy threshold δ , etc., during our discussion, as we focus on applications of diffusion models to high-dimensional data distributions, following the standard practice in the literature.

3 Main Results

Inspired by the acceleration achieved by the parallel sampling technique in LMC and ULMC, we aim to accommodate parallel sampling into the theoretical analysis framework of diffusion models. The benefit of the parallel sampling technique in this scenario has been recently confirmed by up to $14\times$ acceleration achieved by the ParaDiGMS algorithm [1] and ParaTAA [92], where several practical compromises are made to mitigate GPU memory constraints and theoretical guarantees are still lacking.

In this section, we will propose **Parallelized Inference Algorithms for Diffusion Models** with both the **SDE** and probability flow **ODE** implementations, namely the **PIADM-SDE** (Algorithm 1) and **PIADM-ODE** (Algorithm 2), and present theoretical guarantees of our algorithms, including the approximate time complexity and space complexity, for both implementations in Section 3.1 and Section 3.2, respectively. Due to the large number of notations used in the presentation, we give an overview of notations in Appendix A.1 for readers' convenience.

3.1 SDE Implementation

We first focus on the approximation, parallelization strategies, and error analysis of diffusion models with the SDE implementation, *i.e.* the forward and backward process (2.3) and its approximation (2.5) with $v = 1$. We will show that PIADM-SDE achieves an $\tilde{O}(\text{poly log } d)$ approximate time complexity with $\tilde{O}(d^2)$ space complexity.

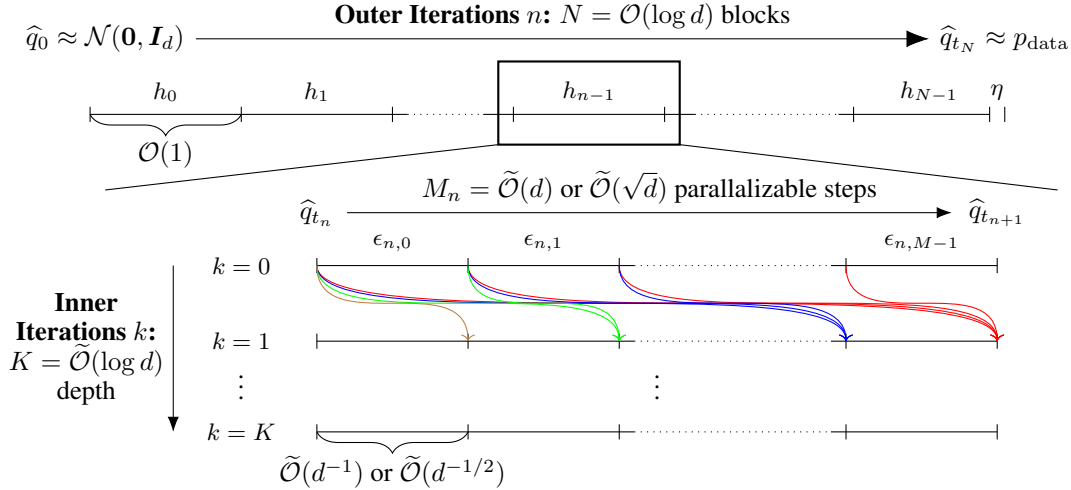


Figure 1: Illustration of PIADM-SDE/ODE. The outer iterations are divided into $\mathcal{O}(\log d)$ blocks of $\mathcal{O}(1)$ length. Within each block, the inner iterations are parallelized with $\tilde{\mathcal{O}}(d)$ steps for SDE (cf. Theorem 3.3), or $\tilde{\mathcal{O}}(\sqrt{d})$ for probability flow ODE implementation (cf. Theorem 3.5). The overall approximate time complexity is $KN = \tilde{\mathcal{O}}(\text{poly log } d)$. brown, green, blue, and red curves represent the computation graph at $t = t_n + \tau_{n,m}$ for $m = 1, 2, M_n - 1, M_n$.

3.1.1 Algorithm

PIADM-SDE is summarized in Algorithm 1 and illustrated in Figure 1. The main idea behind our algorithm is the fact that (2.5) can be efficiently solved by the Picard iteration within a period of $\mathcal{O}(1)$ length, transferring $\tilde{\mathcal{O}}(d)$ sequential computations to a parallelizable iteration of depth $\tilde{\mathcal{O}}(\log d)$. In the following, we introduce the numerical discretization scheme of our algorithm and the implementation of the Picard iteration in detail.

Step Size Scheme. In our algorithm, the time horizon T is first segmented into N blocks of length $(h_n)_{n=0}^{N-1}$, with each $h_n \leq h := T/N = \Omega(1)$, forming a grid $(t_n)_{n=0}^N$ with $t_n = \sum_{j=0}^n h_j$. For any $n \in [0 : N - 1]$, the n -th block is further discretized into a grid $(\tau_{n,m})_{m=0}^{M_n}$ with $\tau_{n,0} = 0$ and $\tau_{n,M_n} = h_n$. We denote the step size of the m -th step in the n -th block as $\epsilon_{n,m} = \tau_{n,m+1} - \tau_{n,m}$, and the total number of steps in the n -th block as M_n .

For the first $N - 1$ blocks, we simply use the unique discretization, i.e. $h_n = h$, $\epsilon_{n,m} = \epsilon$, and $M_n = M := h/\epsilon$, for $n \in [0 : N - 2]$ and $m \in [0 : M - 1]$. Following [104, 107], to curb the potential blow-up of the score function as $t \rightarrow T$, which is shown by [107] for $0 \leq s < t < T$ to be of the order

$$\mathbb{E} \left[\int_s^t \|\nabla \log \tilde{p}_\tau(\tilde{\mathbf{x}}_\tau) - \nabla \log \tilde{p}_s(\tilde{\mathbf{x}}_s)\|^2 d\tau \right] \lesssim d \left(\frac{t-s}{T-t} \right)^2,$$

we apply early stopping at time $t_N = T - \eta$, where η is chosen in a way such that the $\mathcal{O}(\sqrt{\eta})$ 2-Wasserstein distance between \tilde{p}_T and its smoothed version $\tilde{p}_{T-\eta}$ that we aim to sample from alternatively, is tolerable for the downstream tasks. We also impose the exponential decay of the step size towards the data end in the last block. To be specific, we let $h_{N-1} = h - \delta$, and discretize the interval $[t_{N-1}, t_N] = [(N-1)h, T - \eta]$ into a grid $(\tau_{N-1,m})_{m=0}^{M_{N-1}}$ with step sizes $(\epsilon_{N-1,m})_{m=0}^{M_{N-1}-1}$ satisfying

$$\epsilon_{N-1,m} \leq \epsilon \wedge (h - \tau_{N-1,m+1}). \quad (3.1)$$

As shown in Lemma B.7, this exponentially decaying step size scheme towards the data end is crucial to bound the discretization error in the last block.

For the simplicity of notations, we introduce the following indexing function: for $\tau \in [t_n, t_{n+1}]$, we define $I_n(\tau)$ to be the unique integer such that $\sum_{j=1}^{I_n(\tau)} \epsilon_{n,j} \leq \tau < \sum_{j=1}^{I_n(\tau)+1} \epsilon_{n,j}$. We also define

Algorithm 1: PIADM-SDE

Input: $\hat{\mathbf{y}}_0 \sim \hat{q}_0 = \mathcal{N}(0, \mathbf{I}_d)$, a discretization scheme $(T, (h_n)_{n=1}^N$ and $(\tau_{n,m})_{n \in [1:N], m \in [0:M]}$ satisfying (3.1), the depth of iteration K , the learned NN-based score function $\mathbf{s}_t^\theta(\cdot)$.

Output: A sample $\hat{\mathbf{y}}_{t_N} \sim \hat{q}_{t_N} \approx \tilde{p}_T$.

```
1 for  $n = 0$  to  $N - 1$  do
2    $\hat{\mathbf{y}}_{t_n, \tau_{n,m}}^{(0)} \leftarrow \hat{\mathbf{y}}_{t_n}, \boldsymbol{\xi}_m \sim \mathcal{N}(0, \mathbf{I}_d)$  for  $m \in [0 : M_n]$  in parallel;
3   for  $k = 0$  to  $K - 1$  do
4      $\hat{\mathbf{y}}_{t_n, 0}^{(k)} \leftarrow \hat{\mathbf{y}}_{t_n}$ ;
5     for  $m = 0$  to  $M_n$  in parallel do
6        $\hat{\mathbf{y}}_{t_n, \tau_{n,m}}^{(k+1)} \leftarrow e^{\frac{\tau_{n,m}}{2}} \hat{\mathbf{y}}_{t_n, 0}^{(k)}$ 
7          $+ \sum_{j=0}^{m-1} e^{\frac{\tau_{n,m} - \tau_{n,j+1}}{2}} \left[ 2(e^{\epsilon_{n,j}} - 1) \mathbf{s}_{t_n + \tau_{n,j}}^\theta(\hat{\mathbf{y}}_{t_n, \tau_{n,j}}^{(k)}) + \sqrt{e^{\epsilon_{n,j}} - 1} \boldsymbol{\xi}_j \right];$  (3.4)
8     end
9    $\hat{\mathbf{y}}_{t_{n+1}} \leftarrow \hat{\mathbf{y}}_{t_n, \tau_{n, M_n}}^{(K)}$ ;
10 end
```

a piecewise function g such that $g_n(\tau) = \sum_{j=1}^{I_n(\tau)} \epsilon_{n,j}$. It is easy to check that under the uniform discretization for $n \in [1 : N - 1]$, we have $I_n(\tau) = \lfloor \tau/\epsilon \rfloor$ and $g_n(\tau) = \lfloor \tau/\epsilon \rfloor \epsilon$.

Exponential Integrator. For each step $\tau \in [t_n + \tau_{n,m}, t_n + \tau_{n,m+1}]$, we use the following exponential integrator scheme [77], as the numerical discretization of the SDE (2.5):

$$\hat{\mathbf{y}}_{t_n, \tau_{n,m+1}} = e^{\epsilon_{n,m}/2} \hat{\mathbf{y}}_{t_n, \tau_{n,m}} + 2 \left(e^{\epsilon_{n,m}/2} - 1 \right) \mathbf{s}_{t_n + \tau_{n,m}}^\theta(\hat{\mathbf{y}}_{t_n + \tau_{n,m}}) + \sqrt{e^{\epsilon_{n,m}} - 1} \boldsymbol{\xi}_m,$$

where $\boldsymbol{\xi} \sim \mathcal{N}(0, \mathbf{I}_d)$. Lemma B.3 shows its equivalence to approximating (2.5) as

$$d\hat{\mathbf{y}}_{t_n, \tau} = \left[\frac{1}{2} \hat{\mathbf{y}}_{t_n, \tau} + \mathbf{s}_{t_n + \tau_{n,m}}^\theta(\hat{\mathbf{y}}_{t_n, \tau_{n,m}}) \right] d\tau + d\mathbf{w}_{t_n + \tau}, \quad \text{for } \tau \in [\tau_{n,m}, \tau_{n,m+1}]. \quad (3.2)$$

Remark 3.1. One could also implement a straightforward Euler-Maruyama scheme instead of the exponential integrator (3.4), where an additional high-order discretization error term would emerge [104, Theorem 1], which we believe would not affect the overall $\tilde{\mathcal{O}}(\text{poly log } d)$ time complexity with parallel sampling.

Picard Iteration. Within each block, we apply Picard iteration of depth K . As shown by Lemma B.3, the discretized scheme (3.4) implements the following iteration for $k \in [0 : K - 1]$:

$$d\hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} = \left[\frac{1}{2} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} + \mathbf{s}_{t_n + g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(k)}) \right] d\tau + d\mathbf{w}_{t_n + \tau}, \quad \text{for } \tau \in [0, h_n]. \quad (3.3)$$

We denote the distribution of $\hat{\mathbf{y}}_{t_n, \tau}^{(K)}$ by $\hat{q}_{t_n + \tau}$. As proved in Lemma B.6, the iteration above would converge to (3.2) in each block exponentially fast, which given a sufficiently accurate learned score estimation \mathbf{s}_t^θ should be close to the true backward SDE (2.3). One should also notice that the Gaussians $\boldsymbol{\xi}_m$ are only sampled once and used for all iterations.

The parallelization for (3.4) in Algorithm 1 should be understood as that for any $k \in [0 : K - 1]$, each $\mathbf{s}_{t_n + \tau_{n,j}}^\theta(\hat{\mathbf{y}}_{t_n, \tau_{n,j}}^{(k)})$ for $j \in [0 : M_n]$ is evaluated in parallel, with subsequent floating-point operations comparably negligible, resulting in the overall $\mathcal{O}(NK)$ approximate time complexity.

3.1.2 Assumptions

Our theoretical analysis will be built on the following mild assumptions on the regularity of the data distribution and the numerical properties of the neural networks:

Assumption 3.1 ($L^2([0, t_N])$ δ -accurate learned score). The learned NN-based score \mathbf{s}_t^θ is δ_2 -accurate in the sense of

$$\mathbb{E}_{\tilde{p}} \left[\sum_{n=0}^{N-1} \sum_{m=0}^{M_n-1} \epsilon_{n,m} \left\| \mathbf{s}_{t_n+\tau_{n,m}}^\theta(\tilde{\mathbf{x}}_{t_n+\tau_{n,m}}) - \nabla \log \tilde{p}_{t_n+\tau_{n,m}}(\tilde{\mathbf{x}}_{t_n+\tau_{n,m}}) \right\|^2 \right] \leq \delta_2^2. \quad (3.5)$$

Assumption 3.2 (Regular and normalized data distribution). The data distribution p_0 has finite second moments and is normalized such that $\text{cov}_{p_0}(\mathbf{x}_0) = \mathbf{I}_d$.

Assumption 3.3 (Bounded and Lipschitz learned NN-based score). The learned NN-based score function \mathbf{s}_t^θ has bounded C^1 norm, i.e. $\|\mathbf{s}_t^\theta(\cdot)\|_{L^\infty([0,T])} \leq M_s$ with Lipschitz constant L_s .

Remark 3.2. Assumption 3.1 and the finite moment assumption in Assumption 3.2 are standard assumptions across previous theoretical works on diffusion models [100, 104, 111], while we adopt the normalization Assumption 3.2 from [107] to simplify true score function-related computations (cf. Lemma A.8). Assumption 3.3 can be easily satisfied by truncation, ensuring computational stability. Notice that the exponential integrator, one actually applies Picard iteration to $e^{-t/2} \mathbf{s}_t^\theta$, a relaxation of Assumption 3.1 might be possible, which is left for future work.

3.1.3 Theoretical Guarantees

The following theorem summarizes our theoretical analysis for PIADM-SDE (Algorithm 1):

Theorem 3.3 (Theoretical Guarantees for PIADM-SDE). Under Assumptions 3.1, 3.2, and 3.3, given the following choices of the order of the parameters

$$T = \mathcal{O}(\log(d\delta^{-2})), \quad h = \Theta(1), \quad N = \mathcal{O}(\log(d\delta^{-2})), \\ \epsilon = \Theta(d^{-1}\delta^2 \log^{-1}(d\delta^{-2})), \quad M = \mathcal{O}(d\delta^{-2} \log(d\delta^{-2})), \quad K = \tilde{\mathcal{O}}(\log(d\delta^{-2})),$$

and let $L_s^2 h_n e^{\frac{7}{2}h_n} \ll 1$, $\delta_2 \lesssim \delta$, $T \lesssim \log \eta^{-1}$, the distribution \hat{q}_{t_N} that PIADM-SDE (Algorithm 1) generates samples from satisfies the following error bound:

$$D_{\text{KL}}(p_\eta \| \hat{q}_{t_N}) \lesssim d e^{-T} + d \epsilon T + \delta_2^2 + d T e^{-K} \lesssim \delta^2,$$

with a total of $KN = \tilde{\mathcal{O}}(\log^2(d\delta^{-2}))$ approximate time complexity and $dM = \tilde{\mathcal{O}}(d^2\delta^{-2})$ space complexity for parallelizable δ -accurate score function computations.

Remark 3.4. We would like to make the following remarks on the result above:

- The acceleration from $\tilde{\mathcal{O}}(d)$ to $\tilde{\mathcal{O}}(\text{poly log } d)$ is at the cost of a trade-off with extra memory cost of $M = \tilde{\mathcal{O}}(d)$ for computing and updating $\{\mathbf{s}_{t_n+\tau_{n,j}}^\theta(\hat{\mathbf{y}}_{t_n,\tau_{n,m}}^{(k)})\}_{m \in [0:M_n]}$ simultaneously during each Picard iteration;
- Compared with log-concave sampling [130], M being of order $\tilde{\mathcal{O}}(d)$ instead of $\tilde{\mathcal{O}}(\sqrt{d})$ therein is partly due to the time independence of the score function $\nabla \log p(\cdot)$ in general sampling tasks. Besides, the scaling $M = \tilde{\mathcal{O}}(d)$ agrees with the current state-of-the-art dependency [107] for the SDE implementation of diffusion models;
- As mentioned above, the scale of the step size ϵ within one block is still confined to $\Theta(1/M) = \tilde{\mathcal{O}}(1/d)$. The block length h , despite being required to be small compared to $1/L_s$, is of order $\Theta(1)$, resulting in only $\Theta(\log d)$ blocks and thus $\tilde{\mathcal{O}}(\text{poly log } d)$ total iterations.

3.1.4 Proof Sketch

The detailed proof of Theorem 3.3 is deferred to Section B. The pipeline of the proof is to (a) first decompose the error $D_{\text{KL}}(\tilde{p}_{t_N} \| \hat{q}_{t_N})$ into blockwise errors using the chain rule of KL divergence; (b) bound the error in each block by invoking Girsanov's theorem; (c) sum up the errors in all blocks.

The key technical challenge lies in Step (b). Different from all previous theoretical works [100, 104, 111], the Picard iteration in our algorithm generates K paths recursively in each block using the learned score \mathbf{s}_t^θ . And therefore the final path $(\hat{\mathbf{y}}_{t_n,\tau}^{(K)})_{\tau \in [0,h_n]}$ depends on all previous paths $(\hat{\mathbf{y}}_{t_n,\tau}^{(k)})_{\tau \in [0,h_n]}$ for $k \in [0 : K-1]$, ruling out a direct change of measure argument from the

naïve application of Girsanov's theorem. To this end, we need a more sophisticated mathematical framework of stochastic processes, as given in Appendix A.2. We define the measurable space (Ω, \mathcal{F}) with filtrations $(\mathcal{F}_t)_{t \geq 0}$ to specify the probability measures on (Ω, \mathcal{F}) of each Wiener process, and resort to one of the most general forms of Girsanov's theorem ([131, Theorem 8.6.6]). For example, in the n -th block, we apply the following change of measure procedure:

1. Let $q|_{\mathcal{F}_{t_n}}$ be the measure where $w_t(\omega)$ is the shared Wiener process in the Picard iteration (3.3) for any $k \in [0 : K - 1]$;
2. Define another process $d\tilde{w}_{t_n+\tau}(\omega) = dw_{t_n+\tau}(\omega) + \delta_{t_n}(\tau, \omega)d\tau$, where

$$\delta_{t_n}(\tau, \omega) := s_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K-1)}(\omega)) - \nabla \log \tilde{p}_{t_n+\tau}(\hat{\mathbf{y}}_{t_n+\tau}^{(K)}(\omega));$$

3. Invoke Girsanov's theorem, which yields that the Radon-Nikodym derivative of the measure $\tilde{p}|_{\mathcal{F}_{t_n}}$ with respect to $q|_{\mathcal{F}_{t_n}}$ satisfies

$$\log \frac{d\tilde{p}|_{\mathcal{F}_{t_n}}}{dq|_{\mathcal{F}_{t_n}}}(\omega) = - \int_0^{h_n} \delta_{t_n}(\tau, \omega)^\top dw_{t_n+\tau}(\omega) - \frac{1}{2} \int_0^{h_n} \|\delta_{t_n}(\tau, \omega)\|^2 d\tau;$$

4. Conclude that $(\tilde{w}_{t_n+\tau})_{\tau \geq 0}$ is a Wiener process under the measure $\tilde{p}|_{\mathcal{F}_{t_n}}$ and thus (3.3) at iteration K satisfies the following SDE:

$$d\hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) = \left[\frac{1}{2} \hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) + \nabla \log \tilde{p}_{t_n+\tau}(\hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega)) \right] d\tau + d\tilde{w}_{t_n+\tau}(\omega),$$

i.e. the true backward SDE (2.3) with the true score function for $\tau \in [t_n, t_{n+1}]$.

One should notice that this change of measure argument will cause an additional term in the bound of the discrepancy between the first iteration $\hat{\mathbf{y}}_{t_n, \tau}^{(1)}$ and the initial condition $\hat{\mathbf{y}}_{t_n, \tau}^{(0)}$ in Lemma B.5. However, due to the exponential convergence of the Picard iteration, this term does not affect the overall error bound.

3.2 Probability Flow ODE Implementation

In this section, we will show that our parallelization strategy is also compatible with the probability ODE implementation of diffusion models, *i.e.* the forward and backward process (2.3) and its approximation (2.5) with $v = 0$. We will demonstrate that PIADM-ODE (Algorithm 2) further improves the space complexity from $\tilde{\mathcal{O}}(d^2)$ to $\tilde{\mathcal{O}}(d^{3/2})$ while maintaining the same $\tilde{\mathcal{O}}(\text{poly log } d)$ approximate time complexity.

3.2.1 Algorithm

Due to the space limit, we refer the readers to Section C.1 and Algorithm 2 for the details of our parallelization of the probability flow ODE formulation of diffusion models. PIADM-ODE keeps the discretization scheme detailed in Section 3.1.1 that divides the time horizon T into N blocks and uses exponential integrators for all updating rules. Notably, PIADM-ODE has the following distinctions compared with PIADM-SDE (Algorithm 1):

- Instead of applying Picard iteration to the backward SDE as in (3.2), we apply Picard iteration to the probability flow ODE as in (C.3) within each block, which does not require sampling i.i.d. Gaussians to simulate a Wiener process;
- The most significant difference is the adoption of an additional *corrector step* [111] after running the probability flow ODE with Picard iteration within one block. During the corrector step, one augments the state space with a Gaussian that represents the initial momentum and then simulates an underdamped Langevin dynamics for $\mathcal{O}(1)$ time with the learned NN-based score function at the time of the block end;
- We then further parallelize the underdamped Langevin dynamics in the corrector step so that it can also be accomplished with $\mathcal{O}(\log d)$ approximate time complexity, as a naïve implementation would result in $\tilde{\mathcal{O}}(\sqrt{d})$ [130], which is incompatible with our desired poly-logarithmic guarantee.

3.2.2 Assumptions

Due to technicalities specific to this implementation, we need first to modify Assumption 3.1 and add assumption on the Lipschitzness of the true score functions $\nabla \log p_t$, which is a common practice in related literature [104, 111]. Recent work on the probability flow ODE implementation [112, 114] also adopts stronger assumptions compared to the SDE implementation.

Assumption 3.1' ($L^\infty([0, t_N])$ δ -accurate learned score). *For any $n \in [0 : N - 1]$ and $m \in [0 : M_n - 1]$, the learned NN-based score $\mathbf{s}_{t_n+\tau_{n,m}}^\theta$ is δ_∞ -accurate in the sense of*

$$\mathbb{E}_{\tilde{p}_{t_n+\tau_{n,m}}} \left[\left\| \mathbf{s}_{t_n+\tau_{n,m}}^\theta(\tilde{\mathbf{x}}_{t_n+\tau_{n,m}}) - \nabla \log \tilde{p}_{t_n+\tau_{n,m}}(\tilde{\mathbf{x}}_{t_n+\tau_{n,m}}) \right\|^2 \right] \leq \delta_\infty^2.$$

Assumption 3.4 (Bounded and Lipschitz true score). *The true score function $\nabla \log p_t$ has bounded C^1 norm, i.e. $\|\nabla \log p_t(\cdot)\|_{L^\infty([0, T])} \leq M_p$ with Lipschitz constant L_p .*

Further relaxations on Assumption 3.4 to time-dependent assumptions accommodating the blow-up to the data end (e.g. [103, Assumption 1.5]) are left for further work.

3.2.3 Theoretical Guarantees

Our results for PIADM-ODE are summarized in the following theorem:

Theorem 3.5 (Theoretical Guarantees for PIADM-ODE). *Under Assumptions 3.1', 3.2, 3.3, and 3.4, given the following choices of the order of the parameters*

$$T = \mathcal{O}(\log(d\delta^{-2})), \quad h = \Theta(1), \quad N = \mathcal{O}(\log(d\delta^{-2})), \\ \epsilon = \Theta\left(d^{-1/2}\delta \log^{-1}(d^{-1/2}\delta^{-1})\right), \quad M = \mathcal{O}(d^{1/2}\delta^{-1} \log(d^{1/2}\delta^{-1})), \quad K = \tilde{\mathcal{O}}(\log(d\delta^{-2})),$$

for the outer iteration and

$$T^\dagger = \mathcal{O}(1) \lesssim L_p^{-1/2} \wedge L_s^{-1/2}, \quad h^\dagger = \Theta(1), \quad N^\dagger = \mathcal{O}(1), \\ \epsilon^\dagger = \Theta(d^{-1/2}\delta), \quad M^\dagger = \mathcal{O}(d^{1/2}\delta^{-1}), \quad K^\dagger = \mathcal{O}(\log(d\delta^{-2})),$$

for the inner iteration during the corrector step, and let $L_s^2 h^2 e^h \vee L_s^2 h^{\dagger 2} e^{h^\dagger} / \gamma \ll 1$, $\delta_\infty \lesssim \delta \log^{-1}(d\delta^{-2})$, and $\gamma \gtrsim L_p^{1/2}$, then the distribution \hat{q}_{t_N} that PIADM-ODE (Algorithm 2) generates samples from satisfies the following error bound:

$$\text{TV}(p_\eta, \hat{q}_{t_N})^2 \lesssim d e^{-T} + d \epsilon^2 T^2 + (T^2 + N^2) \delta_\infty^2 + d N^2 e^{-K} \lesssim \delta^2,$$

with a total of $(K + K^\dagger N^\dagger)N = \tilde{\mathcal{O}}(\log^2(d\delta^{-2}))$ approximate time complexity and $d(M \vee M^\dagger) = \tilde{\mathcal{O}}(d^{3/2}\delta^{-1})$ space complexity for parallelizable δ -accurate score function computations.

The reduction of space complexity by the probability flow ODE implementation is intuitively owing to the fact that the probability flow ODE process is a deterministic process in time rather than a stochastic process as in the SDE implementation, getting rid of the $\mathcal{O}(\epsilon)$ term derived by Itô's symmetry. This allows the discretization error to be bounded with $\mathcal{O}(\epsilon^2)$ instead (cf. Lemma B.7 and C.5).

3.2.4 Proof Sketch

The details of the proof of Theorem 3.5 are provided in Section C. Along with the complexity benefits the deterministic nature of the probability flow ODE may bring, the analysis is technically more involved than that of Theorem 3.3 and requires an intricate interplay between statistical distances. Several major challenges and our corresponding solutions are summarized below:

- The error of the parallelized algorithm within each block may now only be bounded by 2-Wasserstein distance (cf. Theorem C.7) instead of any f -divergence that admits data processing inequality as in the SDE case by Girsanov's theorem. The additional corrector step exactly handles this issue and would intuitively translate 2-Wasserstein proximity to TV distance proximity (cf. Lemma C.18), allowing the decomposition of the overall error into each block;

- For the corrector step, the underdamped Langevin dynamics as a second-order dynamics requires only $\mathcal{O}(\sqrt{d})$ steps to converge, instead of $\mathcal{O}(d)$ steps in its overdamped counterpart. We then adapt the parallelization technique mentioned in Section 2.2 to conclude that it can be accomplished with $\mathcal{O}(\log d)$ approximate time complexity (cf. Theorem C.17). The error caused by the approximation to the true score and numerical discretization within this step is bounded in KL divergence by invoking Girsanov’s theorem (Theorem A.4) as in the proof of Theorem 3.3;
- Different from the SDE case, where the chain rule of KL divergence can easily decouple the initial distribution and the subsequent dynamics, we need several interpolating processes between the implementation and the true backward process in this case. The final guarantee is in TV distance as it connects with the KL divergence via Pinsker’s inequality and admits data processing inequality. We refer the readers to Figure 2 for an overview of the proof pipeline, as well as the notations and intuitions of the auxiliary and interpolating processes appearing in the proof.

4 Discussion and Conclusion

In this work, we have proposed novel parallelization strategies for the inference of diffusion models in both the SDE and probability flow ODE implementations. Our algorithms, namely PIADM-SDE and PIADM-ODE, are meticulously designed and rigorously proved to achieve $\tilde{\mathcal{O}}(\text{poly } \log d)$ approximate time complexity and $\tilde{\mathcal{O}}(d^2)$ and $\tilde{\mathcal{O}}(d^{3/2})$ space complexity, respectively, marking the first inference algorithm of diffusion and probability flow based models with sub-linear approximate time complexity. Our algorithm intuitively divides the time horizon into several $\mathcal{O}(1)$ blocks and applies Picard iteration within each block in parallel, transferring the time complexity into space complexity. Our analysis is built on a sophisticated mathematical framework of stochastic processes and provides deeper insights into the mathematical theory of diffusion models.

Our findings echo and corroborate the recent empirical work [1, 91–94] that parallel sampling techniques significantly accelerate the inference process of diffusion models. Theoretical exploration of the adaptive block window scheme therein presents an interesting future research potential. Possible future work also includes the investigation of how to apply our parallelization framework to other variants of diffusion models, such as the discrete [23, 132–142] and multi-marginal [143] formulations. Although we anticipate implementing diffusion models in parallel may introduce engineering challenges, *e.g.* scalability, hardware compatibility, memory bandwidth, *etc.*, we believe that our theoretical contributions lay a solid foundation that not only supports but also motivates the empirical development of parallel inference algorithms for diffusion models since advancements continue in GPU power and memory efficiency.

Acknowledgments and Disclosure of Funding

LY acknowledges the support of the National Science Foundation under Grant No. DMS-2011699 and DMS-2208163. GMR is supported by a Google Research Scholar Award.

References

- [1] Andy Shih, Suneel Belkhale, Stefano Ermon, Dorsa Sadigh, and Nima Anari. Parallel sampling of diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [2] Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- [3] Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *arXiv preprint arXiv:2209.15571*, 2022.
- [4] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- [5] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.

- [6] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
- [7] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [8] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *Advances in neural information processing systems*, 34: 1415–1428, 2021.
- [9] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- [10] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [11] Linfeng Zhang, Weinan E, and Lei Wang. Monge-ampère flow for generative modeling. *arXiv preprint arXiv:1809.10188*, 2018.
- [12] Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing diffusion paths for controlled image generation. In *International Conference on Machine Learning*, pages 1737–1752. Pmlr, 2023.
- [13] Xinlei Chen, Zhuang Liu, Saining Xie, and Kaiming He. Deconstructing denoising diffusion models for self-supervised learning. *arXiv preprint arXiv:2401.14404*, 2024.
- [14] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022.
- [15] Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. *arXiv preprint arXiv:2401.08740*, 2024.
- [16] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.
- [17] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. Pmlr, 2021.
- [18] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [19] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [20] Yang Song, Liye Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging with score-based generative models. *arXiv preprint arXiv:2111.08005*, 2021.
- [21] Yu Sun, Zihui Wu, Yifan Chen, Berthy T Feng, and Katherine L Bouman. Provable probabilistic imaging using score-based generative priors. *arXiv preprint arXiv:2310.10835*, 2023.
- [22] Xingyu Xu and Yuejie Chi. Provably robust score-based diffusion posterior sampling for plug-and-play image reconstruction. *arXiv preprint arXiv:2403.17042*, 2024.
- [23] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.

- [24] Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *Advances in Neural Information Processing Systems*, 35:4328–4343, 2022.
- [25] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.
- [26] Yujia Huang, Adishree Ghatare, Yuanzhe Liu, Ziniu Hu, Qinsheng Zhang, Chandramouli S Sastry, Siddharth Gururani, Sageev Oore, and Yisong Yue. Symbolic music generation with non-differentiable rule guided diffusion. *arXiv preprint arXiv:2402.14285*, 2024.
- [27] Gautam Mittal, Jesse Engel, Curtis Hawthorne, and Ian Simon. Symbolic music generation with diffusion models. *arXiv preprint arXiv:2103.16091*, 2021.
- [28] Flavio Schneider. Archisound: Audio generation with diffusion. *arXiv preprint arXiv:2301.13267*, 2023.
- [29] Ruihan Yang, Prakhar Srivastava, and Stephan Mandt. Diffusion probabilistic modeling for video generation. *Entropy*, 25(10):1469, 2023.
- [30] Zihao Li, Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Yinyu Ye, Minshuo Chen, and Mengdi Wang. Diffusion model for data-driven black-box optimization. *arXiv preprint arXiv:2403.13219*, 2024.
- [31] Chen Xu, Jonghyeok Lee, Xiuyuan Cheng, and Yao Xie. Flow-based distributionally robust optimization. *IEEE Journal on Selected Areas in Information Theory*, 2024.
- [32] Ahmed El Alaoui, Andrea Montanari, and Mark Sellke. Sampling from mean-field gibbs measures via diffusion processes. *arXiv preprint arXiv:2310.08912*, 2023.
- [33] Sitan Chen, Vasilis Kontonis, and Kulin Shah. Learning general gaussian mixtures with efficient score matching. *arXiv preprint arXiv:2404.18893*, 2024.
- [34] Ahmed El Alaoui, Andrea Montanari, and Mark Sellke. Sampling from the sherrington-kirkpatrick gibbs measure via algorithmic stochastic localization. In *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 323–334. IEEE, 2022.
- [35] Khashayar Gatmiry, Jonathan Kelner, and Holden Lee. Learning mixtures of gaussians using diffusion models. *arXiv preprint arXiv:2404.18869*, 2024.
- [36] Ye He, Kevin Rojas, and Molei Tao. Zeroth-order sampling methods for non-log-concave distributions: Alleviating metastability by denoising diffusion. *arXiv preprint arXiv:2402.17886*, 2024.
- [37] Xunpeng Huang, Hanze Dong, HAO Yifan, Yian Ma, and Tong Zhang. Reverse diffusion monte carlo. In *The Twelfth International Conference on Learning Representations*, 2023.
- [38] Brice Huang, Andrea Montanari, and Huy Tuan Pham. Sampling from spherical spin glasses in total variation via algorithmic stochastic localization. *arXiv preprint arXiv:2404.15651*, 2024.
- [39] Song Mei and Yuchen Wu. Deep networks as denoising algorithms: Sample-efficient learning of diffusion models in high-dimensional graphical models. *arXiv preprint arXiv:2309.11420*, 2023.
- [40] Andrea Montanari. Sampling, diffusions, and stochastic localization. *arXiv preprint arXiv:2305.10690*, 2023.
- [41] Andrea Montanari and Yuchen Wu. Posterior sampling from the spiked models via diffusion processes. *arXiv preprint arXiv:2304.11449*, 2023.
- [42] Nicholas M Boffi and Eric Vanden-Eijnden. Probability flow solution of the fokker–planck equation. *Machine Learning: Science and Technology*, 4(3):035012, 2023.

- [43] Yan Huang and Li Wang. A score-based particle method for homogeneous landau equation. *arXiv preprint arXiv:2405.05187*, 2024.
- [44] Lingxiao Li, Samuel Hurault, and Justin M Solomon. Self-consistent velocity matching of probability flows. *Advances in Neural Information Processing Systems*, 36, 2024.
- [45] Jianfeng Lu, Yue Wu, and Yang Xiang. Score-based transport modeling for mean-field fokker-planck equations. *Journal of Computational Physics*, 503:112859, 2024.
- [46] Dimitra Maoutsa, Sebastian Reich, and Manfred Opper. Interacting particle solutions of fokker-planck equations through gradient-log-density estimation. *Entropy*, 22(8):802, 2020.
- [47] Amira Alakhdar, Barnabas Poczos, and Newell Washburn. Diffusion models in de novo drug design. *Journal of Chemical Information and Modeling*, 2024.
- [48] Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design. *arXiv preprint arXiv:2402.04997*, 2024.
- [49] Hannes Stark, Bowen Jing, Chenyu Wang, Gabriele Corso, Bonnie Berger, Regina Barzilay, and Tommi Jaakkola. Dirichlet flow matching with applications to dna sequence design. *arXiv preprint arXiv:2402.05841*, 2024.
- [50] Pavel Avdeyev, Chenlai Shi, Yuhao Tan, Kseniia Dudnyk, and Jian Zhou. Dirichlet diffusion score model for biological sequence generation. In *International Conference on Machine Learning*, pages 1276–1301. PMLR, 2023.
- [51] Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Alex X Lu, Nicolo Fusi, Ava P Amini, and Kevin K Yang. Protein generation with evolutionary diffusion: sequence is all you need. *BioRxiv*, pages 2023–09, 2023.
- [52] Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023.
- [53] Jordan Cotler and Semon Rezchikov. Renormalizing diffusion models. *arXiv preprint arXiv:2308.12355*, 2023.
- [54] Florian Furrutter, Gorka Muñoz-Gil, and Hans J Briegel. Quantum circuit synthesis with diffusion models. *arXiv preprint arXiv:2311.02041*, 2023.
- [55] Zhiye Guo, Jian Liu, Yanli Wang, Mengrui Chen, Duolin Wang, Dong Xu, and Jianlin Cheng. Diffusion models in bioinformatics and computational biology. *Nature reviews bioengineering*, 2(2):136–154, 2024.
- [56] Artan Sheshmani, Yi-Zhuang You, Baturalp Buyukates, Amir Ziashahabi, and Salman Avestimehr. Renormalization group flow, optimal transport and diffusion-based generative model. *arXiv preprint arXiv:2402.17090*, 2024.
- [57] Luke Triplett and Jianfeng Lu. Diffusion methods for generating transition paths. *arXiv preprint arXiv:2309.10276*, 2023.
- [58] Lingxiao Wang, Gert Aarts, and Kai Zhou. Generative diffusion models for lattice field theory. *arXiv preprint arXiv:2311.03578*, 2023.
- [59] Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric diffusion model for molecular conformation generation. *arXiv preprint arXiv:2203.02923*, 2022.
- [60] Mengjiao Yang, KwangHwan Cho, Amil Merchant, Pieter Abbeel, Dale Schuurmans, Igor Mordatch, and Ekin Dogus Cubuk. Scalable diffusion for materials generation. *arXiv preprint arXiv:2311.09235*, 2023.

- [61] Jiahao Fan, Ziyao Li, Eric Alcaide, Guolin Ke, Huaqing Huang, and Weinan E. Accurate conformation sampling via protein structural diffusion. *Journal of Chemical Information and Modeling*, 2024.
- [62] Jason Yim, Hannes Stärk, Gabriele Corso, Bowen Jing, Regina Barzilay, and Tommi S Jaakkola. Diffusion models in protein structure and docking. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 14(2):e1711, 2024.
- [63] Yuchen Zhu, Tianrong Chen, Evangelos A Theodorou, Xie Chen, and Molei Tao. Quantum state generation with structure-preserving diffusion model. *arXiv preprint arXiv:2404.06336*, 2024.
- [64] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39, 2023.
- [65] Stanley H Chan. Tutorial on diffusion models for imaging and vision. *arXiv preprint arXiv:2403.18103*, 2024.
- [66] Minshuo Chen, Song Mei, Jianqing Fan, and Mengdi Wang. Opportunities and challenges of diffusion models for generative ai. *National Science Review*, page nwae348, 2024.
- [67] Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34:17695–17709, 2021.
- [68] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [69] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Genie: Higher-order denoising diffusion solvers. *Advances in Neural Information Processing Systems*, 35:30150–30166, 2022.
- [70] Gen Li, Yu Huang, Timofey Efimov, Yuting Wei, Yuejie Chi, and Yuxin Chen. Accelerating convergence of score-based diffusion models, provably. *arXiv preprint arXiv:2403.03852*, 2024.
- [71] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*, 2022.
- [72] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022.
- [73] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022.
- [74] Kaiwen Zheng, Cheng Lu, Jianfei Chen, and Jun Zhu. Dpm-solver-v3: Improved diffusion ode solver with empirical model statistics. *Advances in Neural Information Processing Systems*, 36:55502–55542, 2023.
- [75] Zhenyu Zhou, Defang Chen, Can Wang, and Chun Chen. Fast ode-based sampling for diffusion models in around 5 steps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7777–7786, 2024.
- [76] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35: 26565–26577, 2022.
- [77] Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator. *arXiv preprint arXiv:2204.13902*, 2022.
- [78] Saravanan Kandasamy and Dheeraj Nagaraj. The poisson midpoint method for langevin dynamics: Provably efficient discretization for diffusion models. *arXiv preprint arXiv:2405.17068*, 2024.

- [79] Yuchen Wu, Yuxin Chen, and Yuting Wei. Stochastic runge-kutta methods: Provable acceleration of diffusion models. *arXiv preprint arXiv:2410.04760*, 2024.
- [80] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Score-based generative modeling with critically-damped langevin diffusion. *arXiv preprint arXiv:2112.07068*, 2021.
- [81] Alexia Jolicœur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas. Gotta go fast when generating data with score-based models. *arXiv preprint arXiv:2105.14080*, 2021.
- [82] Hongkai Zheng, Weili Nie, Arash Vahdat, Kamyar Azizzadenesheli, and Anima Anandkumar. Fast sampling of diffusion models via operator learning. In *International Conference on Machine Learning*, pages 42390–42402. PMLR, 2023.
- [83] Yilun Xu, Mingyang Deng, Xiang Cheng, Yonglong Tian, Ziming Liu, and Tommi Jaakkola. Restart sampling for improving generative processes. *Advances in Neural Information Processing Systems*, 36:76806–76838, 2023.
- [84] Jonathan Heek, Emiel Hooeboom, and Tim Salimans. Multistep consistency models. *arXiv preprint arXiv:2403.06807*, 2024.
- [85] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.
- [86] Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models. *arXiv preprint arXiv:2310.14189*, 2023.
- [87] Cheng Lu and Yang Song. Simplifying, stabilizing and scaling continuous-time consistency models. *arXiv preprint arXiv:2410.11081*, 2024.
- [88] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021.
- [89] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14297–14306, 2023.
- [90] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022.
- [91] Hyungjin Chung, Jeongsol Kim, Sehui Kim, and Jong Chul Ye. Parallel diffusion models of operator and image for blind inverse problems. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6059–6069, 2023.
- [92] Zhiwei Tang, Jiasheng Tang, Hao Luo, Fan Wang, and Tsung-Hui Chang. Accelerating parallel sampling of diffusion models. *arXiv preprint arXiv:2402.09970*, 2024.
- [93] Jiezhong Cao, Yue Shi, Kai Zhang, Yulun Zhang, Radu Timofte, and Luc Van Gool. Deep equilibrium diffusion restoration with parallel sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2824–2834, 2024.
- [94] Nikil Rooshan Selvam, Amil Merchant, and Stefano Ermon. Self-refining diffusion samplers: Enabling parallelization via parareal iterations. *arXiv preprint arXiv:2412.08292*, 2024.
- [95] Shivam Gupta, Linda Cai, and Sitan Chen. Faster diffusion-based sampling with randomized midpoints: Sequential and parallel. *arXiv preprint arXiv:2406.00924*, 2024.
- [96] Gen Li and Yuchen Jiao. Improved convergence rate for diffusion probabilistic models. *arXiv preprint arXiv:2410.13738*, 2024.
- [97] Ruoyi Shen and Yin Tat Lee. The randomized midpoint method for log-concave sampling. *Advances in Neural Information Processing Systems*, 32, 2019.

- [98] Belinda Tzen and Maxim Raginsky. Theoretical guarantees for sampling and inference in generative models with latent diffusions. In *Conference on Learning Theory*, pages 3084–3114. PMLR, 2019.
- [99] Adam Block, Youssef Mroueh, and Alexander Rakhlin. Generative modeling with denoising auto-encoders and langevin sampling. *arXiv preprint arXiv:2002.00107*, 2020.
- [100] Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. *arXiv preprint arXiv:2209.11215*, 2022.
- [101] Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence for score-based generative modeling with polynomial complexity. *Advances in Neural Information Processing Systems*, 35:22870–22882, 2022.
- [102] Francesco Pedrotti, Jan Maas, and Marco Mondelli. Improved convergence of score-based diffusion models via prediction-correction. *arXiv preprint arXiv:2305.14164*, 2023.
- [103] Sitan Chen, Giannis Daras, and Alex Dimakis. Restoration-degradation beyond linear diffusions: A non-asymptotic analysis for ddim-type samplers. In *International Conference on Machine Learning*, pages 4462–4484. PMLR, 2023.
- [104] Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *International Conference on Machine Learning*, pages 4735–4763. PMLR, 2023.
- [105] Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for general data distributions. In *International Conference on Algorithmic Learning Theory*, pages 946–985. PMLR, 2023.
- [106] Sokhna Diarra Mbacke and Omar Rivasplata. A note on the convergence of denoising diffusion probabilistic models. *arXiv preprint arXiv:2312.05989*, 2023.
- [107] Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Linear convergence bounds for diffusion models via stochastic localization. *arXiv preprint arXiv:2308.03686*, 2023.
- [108] Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Towards faster non-asymptotic convergence for diffusion-based generative models. *arXiv preprint arXiv:2306.09251*, 2023.
- [109] Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Towards non-asymptotic convergence for diffusion-based generative models. In *The Twelfth International Conference on Learning Representations*, 2024.
- [110] Gen Li, Yuting Wei, Yuejie Chi, and Yuxin Chen. A sharp convergence theory for the probability flow odes of diffusion models. *arXiv preprint arXiv:2408.02320*, 2024.
- [111] Sitan Chen, Sinho Chewi, Holden Lee, Yuanzhi Li, Jianfeng Lu, and Adil Salim. The probability flow ode is provably fast. *Advances in Neural Information Processing Systems*, 36, 2024.
- [112] Xuefeng Gao and Lingjiong Zhu. Convergence analysis for general probability flow odes of diffusion models in wasserstein distances. *arXiv preprint arXiv:2401.17958*, 2024.
- [113] Yuchen Liang, Peizhong Ju, Yingbin Liang, and Ness Shroff. Non-asymptotic convergence of discrete-time diffusion models: New approach and improved rate. *arXiv preprint arXiv:2402.13901*, 2024.
- [114] Daniel Zhengyu Huang, Jiaoyang Huang, and Zhengjiang Lin. Convergence analysis of probability flow ode for score-based generative models. *arXiv preprint arXiv:2404.09730*, 2024.
- [115] Gen Li and Yuling Yan. $o(d/t)$ convergence theory for diffusion probabilistic models under minimal assumptions. *arXiv preprint arXiv:2409.18959*, 2024.

- [116] Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- [117] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- [118] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- [119] Yu Cao, Jingrun Chen, Yixin Luo, and Xiang Zhou. Exploring the optimal choice for generative processes in diffusion models: Ordinary vs stochastic differential equations. *Advances in Neural Information Processing Systems*, 36, 2024.
- [120] Teo Deveney, Jan Stanczuk, Lisa Maria Kreusser, Chris Budd, and Carola-Bibiane Schönlieb. Closing the ode-sde gap in score-based diffusion models through the fokker-planck equation. *arXiv preprint arXiv:2311.15996*, 2023.
- [121] Shen Nie, Hanzhong Allan Guo, Cheng Lu, Yuhao Zhou, Chenyu Zheng, and Chongxuan Li. The blessing of randomness: Sde beats ode in general diffusion-based image editing. *arXiv preprint arXiv:2311.01410*, 2023.
- [122] Charles J Geyer. Markov chain monte carlo maximum likelihood. In *E. M. Kernamidas, editor, Computing Science and Statistics: Proc. 23rd Symposium on the Interface*, pages 156–163. Interface Foundation of North America, 1991.
- [123] Koji Hukushima and Koji Nemoto. Exchange monte carlo method and application to spin glass simulations. *Journal of the Physical Society of Japan*, 65(6):1604–1608, 1996.
- [124] Faming Liang. Use of sequential structure in simulation from high-dimensional systems. *Physical Review E*, 67(5):056101, 2003.
- [125] Nima Anari, Yizhi Huang, Tianyu Liu, Thuy-Duong Vuong, Brian Xu, and Katherine Yu. Parallel discrete sampling via continuous walks. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, pages 103–116, 2023.
- [126] Holden Lee. Parallelising glauher dynamics. *arXiv preprint arXiv:2307.07131*, 2023.
- [127] Lu Yu and Arnak Dalalyana. Parallelized midpoint randomization for langevin monte carlo. *arXiv preprint arXiv:2402.14434*, 2024.
- [128] Ernest Lindelöf. Sur l'application de la méthode des approximations successives aux équations différentielles ordinaires du premier ordre. *Comptes rendus hebdomadaires des séances de l'Académie des sciences*, 116(3):454–457, 1894.
- [129] Emile Picard. Sur les méthodes d'approximations successives dans la théorie des équations différentielles. *American Journal of Mathematics*, pages 87–100, 1898.
- [130] Nima Anari, Sinho Chewi, and Thuy-Duong Vuong. Fast parallel sampling under isoperimetry. *arXiv preprint arXiv:2401.09016*, 2024.
- [131] Bernt Oksendal. *Stochastic differential equations: an introduction with applications*. Springer Science & Business Media, 2013.
- [132] Emiel Hoogeboom, Alexey A Gritsenko, Jasmijn Bastings, Ben Poole, Rianne van den Berg, and Tim Salimans. Autoregressive diffusion models. *arXiv preprint arXiv:2110.02037*, 2021.
- [133] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465, 2021.
- [134] Chenlin Meng, Kristy Choi, Jiaming Song, and Stefano Ermon. Concrete score matching: Generalized score matching for discrete data. *Advances in Neural Information Processing Systems*, 35:34532–34545, 2022.

- [135] Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time discrete diffusion models. *arXiv preprint arXiv:2211.16750*, 2022.
- [136] Pierre H Richemond, Sander Dieleman, and Arnaud Doucet. Categorical sdes with simplex diffusion. *arXiv preprint arXiv:2210.14784*, 2022.
- [137] Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion language modeling by estimating the ratios of the data distribution. *arXiv preprint arXiv:2310.16834*, 2023.
- [138] Griffin Floto, Thorsteinn Jonsson, Mihai Nica, Scott Sanner, and Eric Zhengyu Zhu. Diffusion on the probability simplex. *arXiv preprint arXiv:2309.02530*, 2023.
- [139] Javier E Santos, Zachary R Fox, Nicholas Lubbers, and Yen Ting Lin. Blackout diffusion: generative diffusion models in discrete-state spaces. In *International Conference on Machine Learning*, pages 9034–9059. PMLR, 2023.
- [140] Joe Benton, Yuyang Shi, Valentin De Bortoli, George Deligiannidis, and Arnaud Doucet. From denoising diffusions to denoising markov models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 86(2):286–301, 2024.
- [141] Hongrui Chen and Lexing Ying. Convergence analysis of discrete diffusion model: Exact implementation through uniformization. *arXiv preprint arXiv:2402.08095*, 2024.
- [142] Yinuo Ren, Haoxuan Chen, Grant M Rotskoff, and Lexing Ying. How discrete and continuous diffusion meet: Comprehensive analysis of discrete diffusion models via a stochastic integral framework. *arXiv preprint arXiv:2410.03601*, 2024.
- [143] Michael S Albergo, Nicholas M Boffi, Michael Lindsey, and Eric Vanden-Eijnden. Multi-marginal generative modeling with stochastic interpolants. *arXiv preprint arXiv:2310.03695*, 2023.
- [144] Jean-François Le Gall. *Brownian motion, martingales, and stochastic calculus*. Springer, 2016.
- [145] Arnaud Guillin and Feng-Yu Wang. Degenerate fokker–planck equations: Bismut formula, gradient estimate and harnack inequality. *Journal of Differential Equations*, 253(1):20–40, 2012.

A Mathematical Background

In this section, we will summarize used notations and rigorous mathematical framework of Itô processes as necessary in the proofs. We will also present several technical lemmas for later reference.

A.1 Notations

We adopt the following notations throughout the paper:

Notation	Description
$[a : b]$	The set $\{a, a + 1, \dots, b\}$
\mathbf{I}_d	Identity matrix in $\mathbb{R}^{d \times d}$
$\overleftarrow{*}_t$	$*_{T-t}$
$*$	Used to denote quantities produced by the algorithm
\sim	Used to denote quantities related to the auxiliary processes
$*^\dagger$	Used to denote quantities related to the corrector step
$\ \cdot\ $	The Euclidean norm of a vector
\lesssim or \gtrsim	The inequality holds up to a constant factor
\ll	Absolute continuity (for measures)/ much less than (for quantities)
$(\mathbf{x}_t)_{t \geq 0}$	The forward process of the diffusion model (2.3)
$(\tilde{\mathbf{x}}_t)_{t \in [0, T]}$	The backward process of the diffusion model (2.3)
$(\mathbf{y}_t)_{t \in [0, T]}$	The approximate backward process of the diffusion model (2.5)
$\hat{\mathbf{y}}_{t_n, \tau_n, m}^{(k)}$	The approximate value of the approximate process \mathbf{y}_t at time $t_n + \tau_n, m$ after k iterations in the $(n + 1)$ -th block
$\hat{\mathbf{y}}_{t_n}$	The value of the approximate process \mathbf{y}_t at time t_n
\hat{q}_{t_n}	The distribution of $\hat{\mathbf{y}}_{t_n}$
$\mathbf{z}_{t_1:t_2}$	The path $(\mathbf{z}_t)_{t \in [t_1, t_2]}$ of the process \mathbf{z}_t
$D_f(\cdot \ \cdot)$	The f -divergence between two distributions
$D_{\text{KL}}(\cdot \ \cdot)$	The KL divergence between two distributions
$\text{TV}(\cdot, \cdot)$	The total variation distance between two distributions
$W_2(\cdot, \cdot)$	The 2-Wasserstein distance between two distributions

Table 2: Summary of notations

A.2 Preliminaries

Theorem A.1 (Properties of f -divergence). *Suppose p and q are two probability measures on a common measurable space (Ω, \mathcal{F}) with $p \ll q$. The f -divergence between p and q is defined as*

$$D_f(p \| q) = \mathbb{E}_{X \sim q} \left[f \left(\frac{dp}{dq} \right) \right], \quad (\text{A.1})$$

where $\frac{dp}{dq}$ is the Radon-Nikodym derivative of p with respect to q , and $f : \mathbb{R}^+ \rightarrow \mathbb{R}$ is a convex function. In particular, $D_f(\cdot \| \cdot)$ coincides with the Kullback-Leibler (KL) divergence when $f(x) = x \log x$ and $D_f(\cdot \| \cdot) = \text{TV}$ coincides with the total variation (TV) distance when $f(x) = \frac{1}{2}|x - 1|$.

For the f -divergence defined above, we have the following properties:

1. (Data-processing inequality). Suppose \mathcal{H} is a sub- σ -algebra of \mathcal{F} , the following inequality holds

$$D_f(p|_{\mathcal{H}} \| q|_{\mathcal{H}}) \leq D_f(p \| q);$$

for any f -divergence $D_f(\cdot \| \cdot)$.

2. (Chain rule). Suppose X is a random variable generating a sub- σ -algebra \mathcal{F}_X of \mathcal{F} , and $p(\cdot | X) \ll q(\cdot | X)$ holds for any value of X , then

$$D_{\text{KL}}(p \| q) = D_{\text{KL}}(p|_{\mathcal{F}_X} \| q|_{\mathcal{F}_X}) + \mathbb{E}_{\mathcal{F}_X} [D_{\text{KL}}(p(\cdot | X) \| q(\cdot | X))].$$

In this paper, we consider a probability space (Ω, \mathcal{F}, p) on which $(\mathbf{w}_t(\omega))_{t \geq 0}$ is a Wiener process in \mathbb{R}^d . The Wiener process $(\mathbf{w}_t(\omega))_{t \geq 0}$ generates the filtration $\{\mathcal{F}_t\}_{t \geq 0}$ on the measurable space (Ω, \mathcal{F}) . For an Itô process $\mathbf{z}_t(\omega)$ with the following governing SDE:

$$d\mathbf{z}_t(\omega) = \boldsymbol{\alpha}(t, \omega)dt + \boldsymbol{\Sigma}(t, \omega)d\mathbf{w}_t(\omega),$$

for any time t , we denote the marginal distribution of \mathbf{z}_t by p_t , i.e.

$$p_t := p(\mathbf{z}_t^{-1}(\cdot)), \text{ where } \mathbf{z}_t : \Omega \rightarrow \mathbb{R}^m, \omega \mapsto \mathbf{z}_t(\omega),$$

as well as the *path measure* of the process \mathbf{z}_t in the sense of

$$p_{t_1:t_2} := p(\mathbf{z}_{t_1:t_2}^{-1}(\cdot)), \text{ where } \mathbf{z}_{t_1:t_2} : \Omega \rightarrow \mathcal{C}([t_1, t_2], \mathbb{R}^m), \omega \mapsto (\mathbf{z}_t(\omega))_{t \in [t_1, t_2]}.$$

For the sake of simplicity, we define the following class of functions:

Definition A.2. For any $0 \leq t_1 < t_2$, we define $\mathcal{V}(t_1, t_2)$ as the class of functions $f(t, \omega) : [0, +\infty) \times \Omega \rightarrow \mathbb{R}$ such that

1. $f(t, \omega)$ is $\mathcal{B} \times \mathcal{F}$ -measurable, where \mathcal{B} is the Borel σ -algebra on \mathbb{R}^d ;
2. $f(t, \omega)$ is \mathcal{F}_t -adapted for all $t \geq 0$;
3. The following Novikov condition holds

$$\mathbb{E} \left[\exp \int_{t_1}^{t_2} f^2(t, \omega) dt \right] < +\infty,$$

and $\mathcal{V} = \cap_{t \geq 0} \mathcal{V}(0, t)$. For vectors and matrices, we say it belongs to $\mathcal{V}^n(t, \omega)$ or $\mathcal{V}^{m \times n}(t, \omega)$ if each component of the vector or each entry of the matrix belongs to $\mathcal{V}(t, \omega)$.

Remark A.3. Novikov's condition appeared in the third requirement is often relaxed to the squared integrability condition in the general definition of Itô processes, which requires

$$\mathbb{E} \left[\int_{t_1}^{t_2} f^2(t, \omega) dt \right] < +\infty.$$

Here, we adopt the more restricted condition in the spirit of its necessity for Girsanov's theorem to hold, as we shall see later.

Similar to previous work [111], here we can avoid checking Novikov's condition throughout our proofs below by using the approximation argument presented in [100]. A review of Girsanov can be found in textbooks like in [131, 144]. We will present the following generalized version of Girsanov's theorem:

Theorem A.4 (Girsanov's Theorem [131, Theorem 8.6.6]). Let $\boldsymbol{\alpha}(t, \omega) \in \mathcal{V}^m$, $\boldsymbol{\Sigma}(t, \omega) \in \mathcal{V}^{m \times n}$, and $(\mathbf{w}_t(\omega))_{t \geq 0}$ be a Wiener process on the probability space (Ω, \mathcal{F}, q) . For $t \in [0, T]$, suppose $\mathbf{z}_t(\omega)$ is an Itô process with the following SDE:

$$d\mathbf{z}_t(\omega) = \boldsymbol{\alpha}(t, \omega)dt + \boldsymbol{\Sigma}(t, \omega)d\mathbf{w}_t(\omega), \tag{A.2}$$

and there exist processes $\boldsymbol{\delta}(t, \omega) \in \mathcal{V}^n$ and $\boldsymbol{\beta}(t, \omega) \in \mathcal{V}^m$ such that

1. $\boldsymbol{\Sigma}(t, \omega)\boldsymbol{\delta}(t, \omega) = \boldsymbol{\alpha}(t, \omega) - \boldsymbol{\beta}(t, \omega)$;
2. The process $M_t(\omega)$ as defined below is a martingale with respect to the filtration $\{\mathcal{F}_t\}_{t \geq 0}$ and probability measure q :

$$M_t(\omega) = \exp \left(- \int_0^t \boldsymbol{\delta}(s, \omega)^\top d\mathbf{w}_s(\omega) - \frac{1}{2} \int_0^t \|\boldsymbol{\delta}(s, \omega)\|^2 ds \right),$$

then there exists another probability measure p on (Ω, \mathcal{F}) such that

1. $p \ll q$ with the Radon-Nikodym derivative $\frac{dp}{dq}(\omega) = M_T(\omega)$,

2. The process $\tilde{\mathbf{w}}_t(\omega)$ as defined below is a Wiener process on (Ω, \mathcal{F}, p) :

$$\tilde{\mathbf{w}}_t(\omega) = \mathbf{w}_t(\omega) + \int_0^t \boldsymbol{\delta}(s, \omega) ds,$$

3. Any continuous path in $\mathcal{C}([t_1, t_2], \mathbb{R}^m)$ generated by the process \mathbf{z}_t satisfies the following SDE under the probability measure p :

$$d\tilde{\mathbf{z}}_t(\omega) = \boldsymbol{\beta}(t, \omega)dt + \boldsymbol{\Sigma}(t, \omega)d\tilde{\mathbf{w}}_t(\omega). \quad (\text{A.3})$$

Corollary A.5. Suppose the conditions in Theorem A.4 hold, then for any $t_1, t_2 \in [0, T]$ with $t_1 < t_2$, the path measure of the SDE (A.3) under the probability measure p in the sense of $p_{t_1:t_2} = p(\mathbf{z}_{t_1:t_2}^{-1}(\cdot))$ is absolutely continuous with respect to the path measure of the SDE (A.2) in the sense of $q_{t_1:t_2} = q(\mathbf{z}_{t_1:t_2}^{-1}(\cdot))$. Moreover, the KL divergence between the two path measures is given by

$$D_{\text{KL}}(p_{t_1:t_2} \| q_{t_1:t_2}) = D_{\text{KL}}(p_{t_1} \| q_{t_1}) + \mathbb{E}_{\omega \sim p|_{\mathcal{F}_{t_1}}} \left[\frac{1}{2} \int_{t_1}^{t_2} \|\boldsymbol{\delta}(t, \omega)\|^2 dt \right] \quad (\text{A.4})$$

Proof. First, by Theorem A.1, we have

$$D_{\text{KL}}(p_{t_1:t_2} \| q_{t_1:t_2}) = D_{\text{KL}}(p|_{\mathcal{F}_{t_1}} \| q|_{\mathcal{F}_{t_1}}) + \mathbb{E}_{\mathbf{z} \sim p|_{\mathcal{F}_{t_1}}} \left[D_{\text{KL}}(p(\tilde{\mathbf{z}}_{t_1:t_2}^{-1}(\cdot)) | \tilde{\mathbf{z}}_{t_1} = \tilde{\mathbf{z}}) \| q(\tilde{\mathbf{z}}_{t_1:t_2}^{-1}(\cdot)) | \tilde{\mathbf{z}}_{t_1} = \tilde{\mathbf{z}}) \right].$$

From Girsanov's theorem (Theorem A.4), we have that the measure $p|_{\mathcal{F}_{t_1}}$ is absolutely continuous with respect to $q|_{\mathcal{F}_{t_1}}$, which allows us to compute the second term above as follows:

$$\begin{aligned} & D_{\text{KL}}(p(\tilde{\mathbf{z}}_{t_1:t_2}^{-1}(\cdot) | \tilde{\mathbf{z}}_{t_1} = \tilde{\mathbf{z}}) \| q(\tilde{\mathbf{z}}_{t_1:t_2}^{-1}(\cdot) | \tilde{\mathbf{z}}_{t_1} = \tilde{\mathbf{z}})) \\ &= \mathbb{E}_{\tilde{\mathbf{z}}_{t_1:t_2}} \left[\log \frac{dp(\tilde{\mathbf{z}}_{t_1:t_2}^{-1}(\cdot) | \tilde{\mathbf{z}}_{t_1} = \mathbf{z})}{dq(\tilde{\mathbf{z}}_{t_1:t_2}^{-1}(\cdot) | \tilde{\mathbf{z}}_{t_1} = \mathbf{z})} \right] = \mathbb{E}_{\omega \sim p|_{\mathcal{F}_{t_1}}} \left[\log \frac{dp|_{\mathcal{F}_{t_1}}}{dq|_{\mathcal{F}_{t_1}}} \right] \\ &= \mathbb{E}_{\omega \sim p|_{\mathcal{F}_{t_1}}} \left[- \int_{t_1}^{t_2} \boldsymbol{\delta}(t, \omega)^\top d\mathbf{w}_t(\omega) - \frac{1}{2} \int_{t_1}^{t_2} \|\boldsymbol{\delta}(t, \omega)\|^2 dt \right] \\ &= \mathbb{E}_{\omega \sim p|_{\mathcal{F}_{t_1}}} \left[- \int_{t_1}^{t_2} \boldsymbol{\delta}(t, \omega)^\top (d\tilde{\mathbf{w}}_t(\omega) - \boldsymbol{\delta}(t, \omega)dt) - \frac{1}{2} \int_{t_1}^{t_2} \|\boldsymbol{\delta}(t, \omega)\|^2 dt \right] \\ &= \mathbb{E}_{\omega \sim p|_{\mathcal{F}_{t_1}}} \left[\frac{1}{2} \int_{t_1}^{t_2} \|\boldsymbol{\delta}(t, \omega)\|^2 dt \right], \end{aligned}$$

and therefore

$$D_{\text{KL}}(p_{t_1:t_2} \| q_{t_1:t_2}) = D_{\text{KL}}(p_{t_1} \| q_{t_1}) + \mathbb{E}_{\omega \sim p|_{\mathcal{F}_{t_1}}} \left[\frac{1}{2} \int_{t_1}^{t_2} \|\boldsymbol{\delta}(t, \omega)\|^2 dt \right],$$

which completes the proof. \square

A.3 Helper Lemmas

Lemma A.6 ([107, Lemma 2]). For the backward process (2.3), we have for $0 \leq s < t < T$,

$$\frac{d}{dt} (\mathbb{E} [\|\nabla \log \tilde{p}_t(\tilde{\mathbf{x}}_t) - \nabla \log \tilde{p}_s(\tilde{\mathbf{x}}_s)\|^2]) \leq \frac{1}{2} \mathbb{E} [\|\nabla \log \tilde{p}_s(\tilde{\mathbf{x}}_s)\|^2] + \mathbb{E} [\|\nabla^2 \log \tilde{p}_t(\tilde{\mathbf{x}}_t)\|_F^2].$$

Lemma A.7 ([107, Lemma 3]). For the forward process (2.3), we have for $0 \leq t < T$,

$$\mathbb{E} [\nabla \log p_t(\mathbf{x}_t)] \leq d\sigma_t^{-2}, \text{ and } \mathbb{E} [\|\nabla^2 \log p_t(\mathbf{x}_t)\|_F^2] \leq d\sigma_t^{-4} + 2 \frac{d}{dt} (\sigma_t^{-4} \mathbb{E} [\text{tr } \boldsymbol{\Sigma}_t]),$$

where the posterior covariance matrix $\boldsymbol{\Sigma}_t := \text{cov}_{p_{0|t}}(\mathbf{x}_0)$ and $\sigma_t^2 = 1 - e^{-t}$. Moreover, the posterior covariance matrix $\boldsymbol{\Sigma}_t$ satisfies

$$\mathbb{E} [\text{tr } \boldsymbol{\Sigma}_t] \lesssim d \wedge d\sigma_t^2.$$

Lemma A.8. For any $n \in [0 : N - 1]$ and $\tau \in [0, h_n]$, under the assumption $\text{cov}_{p_0}(\mathbf{x}_0) = \mathbf{I}_d$, we have

$$\mathbb{E} [\|\tilde{\mathbf{x}}_{t_n}\|^2] \leq 2d, \quad (\text{A.5})$$

and

$$\mathbb{E} [\|\tilde{\mathbf{x}}_{t_n} - \tilde{\mathbf{x}}_{t_n+\tau}\|^2] \leq 3d. \quad (\text{A.6})$$

Proof. Conditioned on \mathbf{x}_0 , we have that

$$\tilde{\mathbf{x}}_{t_n} = \mathbf{x}_{T-t_n} \sim \mathcal{N} \left(e^{-\frac{1}{2}(T-t_n)} \mathbf{x}_0, (1 - e^{-(T-t_n)}) \mathbf{I}_d \right),$$

and

$$\tilde{\mathbf{x}}_{t_n+\tau} = \mathbf{x}_{T-t_n-\tau} \sim \mathcal{N} \left(e^{-\frac{1}{2}(T-t_n-\tau)} \mathbf{x}_0, (1 - e^{-(T-t_n-\tau)}) \mathbf{I}_d \right)$$

for any $\tau \in [0, h_n]$. Therefore, we have

$$\begin{aligned} \mathbb{E} [\|\tilde{\mathbf{x}}_{t_n}\|^2] &= \mathbb{E} [\mathbb{E} [\|\mathbf{x}_{T-t_n}\|^2 | \mathbf{x}_0]] \\ &\leq \mathbb{E} \left[\mathbb{E} [\|\mathbf{x}_{T-t_n} - e^{-\frac{1}{2}(T-t_n)} \mathbf{x}_0\|^2 | \mathbf{x}_0] + \|e^{-\frac{1}{2}(T-t_n)} \mathbf{x}_0\|^2 \right] \\ &\leq d(1 - e^{-(T-t_n)}) + e^{-(T-t_n)} \mathbb{E} [\|\mathbf{x}_0\|^2] \leq 2d. \end{aligned}$$

Taking the difference between them then implies that for any $\tau \in [0, h_n]$,

$$\begin{aligned} \mathbb{E} [\|\tilde{\mathbf{x}}_{t_n} - \tilde{\mathbf{x}}_{t_n+\tau}\|^2] &= \mathbb{E} [\mathbb{E} [\|\mathbf{x}_{T-t_n} - \mathbf{x}_{T-t_n-\tau}\|^2 | \mathbf{x}_0]] \\ &\leq d(2 - e^{-(T-t_n)} - e^{-(T-t_n-\tau)}) \\ &\quad + \left(e^{-\frac{1}{2}(T-t_n)} - e^{-\frac{1}{2}(T-t_n-\tau)} \right)^2 \mathbb{E} [\|\mathbf{x}_0\|^2] \\ &\leq 2d + e^{-(T-t_n-\tau)} (1 - e^{-\frac{1}{2}\tau})^2 \mathbb{E} [\|\mathbf{x}_0\|^2] \leq 3d. \end{aligned}$$

□

Lemma A.9 (Lemma 9 in [104]). For $\hat{q}_0 \sim \mathcal{N}(0, \mathbf{I}_d)$ and $\tilde{p}_0 = p_T$ is the distribution of the solution to the forward process (2.3), we have

$$\text{TV}(\tilde{p}_0, \hat{q}_0)^2 \leq D_{\text{KL}}(\tilde{p}_0 \| \hat{q}_0) \lesssim d e^{-T}.$$

B Details of SDE Implementation

In this section, we will present the missing proofs for Theorem 3.3. For readers' convenience, we reiterate the backward process (2.3)

$$d\tilde{\mathbf{x}}_t = \left[\frac{1}{2} \tilde{\mathbf{x}}_t + \nabla \log \tilde{p}_t(\tilde{\mathbf{x}}_t) \right] dt + d\mathbf{w}_t, \quad \text{with } \tilde{\mathbf{x}}_0 \sim p_T, \quad (\text{B.1})$$

and its approximate version (2.5) with the learned score function

$$d\mathbf{y}_t = \left[\frac{1}{2} \mathbf{y}_t + \mathbf{s}_t^\theta(\mathbf{y}_t) \right] dt + d\mathbf{w}_t, \quad \text{with } \mathbf{y}_0 \sim \mathcal{N}(0, \mathbf{I}_d).$$

The filtration \mathcal{F}_t refers to the filtration of the SDE (B.1) up to time t .

B.1 Auxiliary Process

We would like first to consider the errors that Algorithm 1 may cause within one block of update. To this end, we consider the following auxiliary process for $\tau \in [0, h_n]$ conditioned on the filtration \mathcal{F}_{t_n} at time t_n :

Definition B.1 (Auxiliary Process). *For any $n \in [0 : N - 1]$, we define the auxiliary process $(\hat{\mathbf{y}}_{t_n, \tau}^{(k)})_{\tau \in [0, h_n]}$ as the solution to the following SDE recursively for $k \in [0 : K - 1]$:*

$$d\hat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega) = \left[\frac{1}{2} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega) + \mathbf{s}_{t_n + g_n(\tau)}^\theta \left(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(k)}(\omega) \right) \right] d\tau + d\mathbf{w}_{t_n + \tau}(\omega), \quad (\text{B.2})$$

with the initial condition

$$\hat{\mathbf{y}}_{t_n, \tau}^{(0)}(\omega) \equiv \hat{\mathbf{y}}_{t_n}(\omega) \text{ for } \tau \in [0, h_n], \quad \text{and} \quad \hat{\mathbf{y}}_{t_n, 0}^{(k)}(\omega) \equiv \hat{\mathbf{y}}_{t_n}(\omega) \text{ for } k \in [1 : K] \quad (\text{B.3})$$

where $\hat{\mathbf{y}}_{t_n}(\omega) = \hat{\mathbf{y}}_{t_{n-1}, \tau_{n-1}, M_{n-1}}^{(K)}(\omega)$ if $n \in [1 : N - 1]$ and $\hat{\mathbf{y}}_{t_0}(\omega) \sim \mathcal{N}(0, \mathbf{I}_d)$.

The iteration should be perceived as a deterministic procedure to each event $\omega \in \Omega$, i.e. each realization of the Wiener process $(\mathbf{w}_t)_{t \geq 0}$. The following lemma clarifies this fact and proves the well-definedness and parallelability of the iteration in (B.2).

Lemma B.2. *The auxiliary process $(\hat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega))_{\tau \in [0, h_n]}$ is $\mathcal{F}_{t_n + \tau}$ -adapted for any $k \in [0 : K]$ and $n \in [0 : N - 1]$.*

Proof. Since the initialization $\hat{\mathbf{y}}_{t_n, \tau}^{(0)}(\omega) \equiv \hat{\mathbf{y}}_{t_n}(\omega)$ for $\tau \in [0, h_n]$, where $\hat{\mathbf{y}}_{t_n}(\omega)$ is \mathcal{F}_{t_n} -adapted, it is obvious that $\hat{\mathbf{y}}_{t_n, \tau}^{(0)}(\omega)$ is $\mathcal{F}_{t_n + \tau}$ -adapted. Now suppose that $(\hat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega))_{\tau \in [0, h_n]}$ is $\mathcal{F}_{t_n + \tau}$ -adapted, since $g_n(\tau) \leq \tau$, we have the following Itô integral well-defined and $\mathcal{F}_{t_n + \tau}$ -adapted:

$$\int_0^\tau \mathbf{s}_{t_n + g_n(\tau')}^\theta \left(\hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)}(\omega) \right) d\tau',$$

and therefore (B.2) has a unique strong solution $(\hat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega))_{\tau \in [0, h_n]}$ that is also $\mathcal{F}_{t_n + \tau}$ -adapted. The lemma follows by induction. \square

Lemma B.3 (Equivalence between (3.4) and (B.2)). *For any $n \in [0 : N - 1]$, the update rule (3.4) in Algorithm 1 is equivalent to the exact solution of the auxiliary process (B.2) for any $k \in [0 : K - 1]$ and $\tau \in [0, h_n]$.*

Proof. The dependency on ω will be omitted in the proof below.

Rewriting (B.2) and multiplying $e^{-\frac{\tau}{2}}$ on both sides yield

$$d \left[e^{-\frac{\tau}{2}} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} \right] = e^{-\frac{\tau}{2}} \left[d\hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \frac{1}{2} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} d\tau \right] = e^{-\frac{\tau}{2}} \left[\mathbf{s}_{t_n + g_n(\tau)}^\theta \left(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(k)} \right) d\tau + d\mathbf{w}_{t_n + \tau} \right].$$

Integrating on both sides from 0 to τ implies

$$\begin{aligned} e^{-\frac{\tau}{2}} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \hat{\mathbf{y}}_{t_n, 0}^{(k+1)} &= \int_0^\tau e^{-\frac{\tau'}{2}} \left(\mathbf{s}_{t_n + g_n(\tau')}^\theta \left(\hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)} \right) d\tau' + d\mathbf{w}_{t_n + \tau'} \right) \\ &= \sum_{m=0}^{M_n} \int_{\tau \wedge t_{n, m}}^{\tau \wedge \tau_{n, m+1}} e^{-\frac{\tau'}{2}} \mathbf{s}_{t_n + \tau_{n, m}}^\theta \left(\hat{\mathbf{y}}_{t_n, \tau_{n, m}}^{(k)} \right) d\tau' + \int_0^\tau e^{-\frac{\tau'}{2}} d\mathbf{w}_{t_n + \tau'} \\ &= \sum_{m=0}^{M_n} 2 \left(e^{-\frac{\tau \wedge \tau_{n, m}}{2}} - e^{-\frac{\tau \wedge \tau_{n, m+1}}{2}} \right) \mathbf{s}_{t_n + \tau_{n, m}}^\theta \left(\hat{\mathbf{y}}_{t_n, \tau_{n, m}}^{(k)} \right) + \int_0^\tau e^{-\frac{\tau'}{2}} d\mathbf{w}_{t_n + \tau'}, \end{aligned}$$

and then multiplying $e^{\frac{\tau}{2}}$ on both sides above yields

$$\begin{aligned} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} &= e^{\frac{\tau}{2}} \hat{\mathbf{y}}_{t_n, 0}^{(k+1)} + \sum_{m=0}^{M_n} 2 \left(e^{\frac{\tau \wedge \tau_{n, m+1} - \tau \wedge \tau_{n, m}}{2}} - 1 \right) e^{\frac{0 \vee (\tau - \tau_{n, m+1})}{2}} \mathbf{s}_{t_n + \tau_{n, m}}^\theta \left(\hat{\mathbf{y}}_{t_n, \tau_{n, m}}^{(k)} \right) \\ &\quad + \sum_{m=0}^{M_n} \int_{\tau \wedge t_{n, m}}^{\tau \wedge \tau_{n, m+1}} e^{\frac{\tau - \tau'}{2}} d\mathbf{w}_{t_n + \tau'}, \end{aligned}$$

where, by Itô isometry, we have

$$\int_{\tau \wedge \tau_{n,m}}^{\tau \wedge \tau_{n,m+1}} e^{\frac{\tau - \tau'}{2}} d\mathbf{w}_{t_n + \tau'} \sim \mathcal{N}\left(\mathbf{0}, (e^{\tau \wedge \tau_{n,m+1} - \tau \wedge \tau_{n,m}} - 1) e^{0 \vee (\tau - \tau_{n,m+1})} \mathbf{I}_d\right)$$

for $\tau > \tau_{n,m}$ and equals to $\mathbf{0}$ otherwise. Plugging in $\tau = \tau_{j,m}$ gives us (3.4), as desired. \square

B.2 Errors within Block

We shall invoke Girsanov's theorem (Theorem A.4) in the procedure as detailed below:

1. Setting (A.2) in Theorem A.4 as the auxiliary process (B.2) at iteration K , where $\mathbf{w}_t(\omega)$ is a Wiener process under the measure $q|_{\mathcal{F}_{t_n}}$;
2. Defining another process $\tilde{\mathbf{w}}_{t_n + \tau}(\omega)$ governed by the following SDE:

$$d\tilde{\mathbf{w}}_{t_n + \tau}(\omega) = d\mathbf{w}_{t_n + \tau}(\omega) + \boldsymbol{\delta}_{t_n}(\tau, \omega) d\tau,$$

where

$$\boldsymbol{\delta}_{t_n}(\tau, \omega) := \mathbf{s}_{t_n + g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K-1)}(\omega)) - \nabla \log \tilde{p}_{t_n + \tau}(\hat{\mathbf{y}}_{t_n + \tau}^{(K)}(\omega)), \quad (\text{B.4})$$

and computing the Radon-Nikodym derivative of the measure $\tilde{p}|_{\mathcal{F}_{t_n}}$ with respect to $q|_{\mathcal{F}_{t_n}}$ as

$$\frac{d\tilde{p}|_{\mathcal{F}_{t_n}}}{dq|_{\mathcal{F}_{t_n}}}(\omega) := \exp\left(-\int_0^{h_n} \boldsymbol{\delta}_{t_n}(\tau, \omega)^\top d\mathbf{w}_{t_n + \tau}(\omega) - \frac{1}{2} \int_0^{h_n} \|\boldsymbol{\delta}_{t_n}(\tau, \omega)\|^2 d\tau\right),$$

3. Concluding that (B.2) at iteration K under the measure $q|_{\mathcal{F}_{t_n}}$ satisfies the following SDE:

$$d\hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) = \left[\frac{1}{2} \hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) + \nabla \log \tilde{p}_{t_n + \tau}(\hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega))\right] d\tau + d\tilde{\mathbf{w}}_{t_n + \tau}(\omega),$$

with $(\tilde{\mathbf{w}}_{t_n + \tau})_{\tau \geq 0}$ being a Wiener process under the measure $\tilde{p}|_{\mathcal{F}_{t_n}}$. If we replace $\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)$ by $\tilde{\mathbf{x}}_{t_n + \tau}(\omega)$, one should notice (B.5) is immediately the original backward SDE (2.3) with the true score function on $t \in [t_n, t_{n+1}]$:

$$d\tilde{\mathbf{x}}_{t_n + \tau}(\omega) = \left[\frac{1}{2} \tilde{\mathbf{x}}_{t_n + \tau}(\omega) + \nabla \log \tilde{p}_{t_n + \tau}(\tilde{\mathbf{x}}_{t_n + \tau}(\omega))\right] d\tau + d\tilde{\mathbf{w}}_{t_n + \tau}(\omega). \quad (\text{B.5})$$

Remark B.4. The applicability of Girsanov's theorem here relies on the \mathcal{F}_τ -adaptivity of $\mathbf{s}_{t_n + g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K-1)}(\omega))$ established by Lemma B.2. One should notice the change of measure procedure above depends on the number of iterations K , and different K would lead to different transform (B.4).

Then Corollary A.5 provides the following computation

$$\begin{aligned} D_{\text{KL}}(\tilde{p}_{t_{n+1}} \|\hat{q}_{t_{n+1}}) &\leq D_{\text{KL}}(\tilde{p}_{t_n:t_{n+1}} \|\hat{q}_{t_n:t_{n+1}}) \\ &= D_{\text{KL}}(\tilde{p}_{t_n} \|\hat{q}_{t_n}) + \mathbb{E}_{\omega \sim q|_{\mathcal{F}_{t_n}}} \left[\frac{1}{2} \int_0^{h_n} \|\boldsymbol{\delta}_{t_n}(\tau, \omega)\|^2 d\tau \right], \end{aligned} \quad (\text{B.6})$$

where the first inequality is by the data-processing inequality (Theorem A.1). Now, the problem remaining is to bound the discrepancy quantified by

$$\begin{aligned}
& \int_0^{h_n} \|\delta_{t_n}(\tau, \omega)\|^2 d\tau \\
&= \int_0^{h_n} \left\| \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K-1)}(\omega)) - \nabla \log \tilde{p}_{t_n+\tau}(\hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega)) \right\|^2 d\tau \\
&\leq 3 \underbrace{\left(\int_0^{h_n} \left\| \nabla \log \tilde{p}_{t_n+g_n(\tau)}(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) - \nabla \log \tilde{p}_{t_n+\tau}(\hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega)) \right\|^2 d\tau \right)}_{:= A_{t_n}(\omega)} \\
&\quad + \underbrace{\left(\int_0^{h_n} \left\| \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) - \nabla \log \tilde{p}_{t_n+g_n(\tau)}(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) \right\|^2 d\tau \right)}_{:= B_{t_n}(\omega)} \\
&\quad + \int_0^{h_n} \left\| \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) - \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K-1)}(\omega)) \right\|^2 d\tau.
\end{aligned} \tag{B.7}$$

Before we continue our proof, we would like first to provide the following lemma bounding the behavior of the auxiliary process (B.2) when $k = 0$ for $\tau \in [0, h_n]$.

Lemma B.5. *For any $n \in [0 : N - 1]$, suppose the initialization $\hat{\mathbf{y}}_{t_n}$ in (B.3) of the auxiliary process (B.2) follows the distribution of $\tilde{\mathbf{x}}_{t_n} \sim \tilde{p}_{t_n}$, then the following estimate holds*

$$\begin{aligned}
& \sup_{\tau \in [0, h_n]} \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\|\hat{\mathbf{y}}_{t_n, \tau}^{(1)}(\omega) - \hat{\mathbf{y}}_{t_n, \tau}^{(0)}(\omega)\|^2 \right] \\
&\leq h_n e^{\frac{7}{2} h_n} (M_s^2 + 2d) + 3e^{\frac{7}{2} h_n} \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} [A_{t_n}(\omega) + B_{t_n}(\omega)] \\
&\quad + 3e^{\frac{7}{2} h_n} h_n L_s^2 \sup_{\tau \in [0, h_n]} \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\left\| \hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) - \hat{\mathbf{y}}_{t_n, \tau}^{(K-1)}(\omega) \right\|^2 \right].
\end{aligned} \tag{B.8}$$

Proof. Let $\mathbf{z}_{t_n, \tau} = \hat{\mathbf{y}}_{t_n, \tau}^{(1)} - \hat{\mathbf{y}}_{t_n, \tau}^{(0)}$. For $k = 0$, we can rewrite (B.2) as

$$d\mathbf{z}_{t_n, \tau} = \left[\frac{1}{2} \left(\mathbf{z}_{t_n, \tau} + \hat{\mathbf{y}}_{t_n, \tau}^{(0)} \right) + \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(0)}) \right] d\tau + d\mathbf{w}_{t_n+\tau},$$

By applying Itô's lemma and plugging in the expression of $\mathbf{w}_{t_n+\tau}$ given by Theorem A.4, we have

$$\begin{aligned}
d\|\mathbf{z}_{t_n, \tau}\|^2 &= \left[\|\mathbf{z}_{t_n, \tau}\|^2 + \mathbf{z}_{t_n, \tau}^\top \hat{\mathbf{y}}_{t_n, \tau}^{(0)} + 2\mathbf{z}_{t_n, \tau}^\top \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(0)}) + d \right] d\tau \\
&\quad + 2\mathbf{z}_{t_n, \tau}^\top (d\tilde{\mathbf{w}}_{t_n+\tau}(\omega) - \delta_{t_n}(\tau, \omega) d\tau),
\end{aligned} \tag{B.9}$$

By integrating from 0 to τ and taking the expectation on both sides of (B.9), we obtain that

$$\begin{aligned}
& \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} [\|\mathbf{z}_{t_n, \tau}\|^2] \\
&= \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\int_0^\tau \left(\|\mathbf{z}_{t_n, \tau'}\|^2 + \mathbf{z}_{t_n, \tau'}^\top \hat{\mathbf{y}}_{t_n, \tau'}^{(0)} + 2\mathbf{z}_{t_n, \tau'}^\top \mathbf{s}_{t_n+g_n(\tau')}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(0)}) + d \right) d\tau' \right] \\
&\quad + 2\mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\int_0^\tau \mathbf{z}_{t_n, \tau'}^\top (d\tilde{\mathbf{w}}_{t_n+\tau'}(\omega) - \delta_{t_n}(\tau', \omega) d\tau') \right],
\end{aligned}$$

and by AM-GM, we further have

$$\begin{aligned}
& \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} [\|\mathbf{z}_{t_n, \tau}\|^2] \\
&\leq \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\int_0^\tau \left[\frac{7}{2} \|\mathbf{z}_{t_n, \tau'}\|^2 + \frac{1}{2} \|\hat{\mathbf{y}}_{t_n, \tau'}^{(0)}\|^2 + \left\| \mathbf{s}_{t_n+g_n(\tau')}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(0)}) \right\|^2 + d + \|\delta_{t_n}(\tau, \omega)\|^2 \right] d\tau' \right] \\
&\leq \int_0^\tau \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\frac{7}{2} \|\mathbf{z}_{t_n, \tau'}\|^2 + \|\delta_{t_n}(\tau, \omega)\|^2 \right] d\tau' + \left(\frac{1}{2} \mathbb{E} \left[\|\hat{\mathbf{y}}_{t_n, \tau}^{(0)}\|^2 \right] + M_s^2 + d \right) \tau,
\end{aligned}$$

where $\delta_{t_n}(\tau, \omega)$ is defined in (B.4). Similar to (B.7), we may use triangle inequality to upper bound $\|\delta_{t_n}(\tau, \omega)\|^2$, which implies that for any $\tau \in [0, h_n]$

$$\begin{aligned} & \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [\|\mathbf{z}_{t_n, \tau}\|^2] \\ & \leq \frac{7}{2} \int_0^\tau \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [\|\mathbf{z}_{t_n, \tau'}\|^2] d\tau' + \left(\frac{1}{2} \mathbb{E} [\|\hat{\mathbf{y}}_{t_n, \tau}^{(0)}\|^2] + M_s^2 + d \right) \tau \\ & + 3 \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\int_0^\tau \left\| \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K-1)}(\omega)) - \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) \right\|^2 d\tau' \right] \\ & + 3 \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\int_0^\tau \left\| \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) - \nabla \log \tilde{p}_{t_n, g_n(\tau)}(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) \right\|^2 d\tau' \right] \\ & + 3 \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\int_0^\tau \left\| \nabla \log \tilde{p}_{t_n+g_n(\tau)}(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) - \nabla \log \tilde{p}_{t_n+\tau}(\hat{\mathbf{y}}_{t_n+\tau}^{(K)}(\omega)) \right\|^2 d\tau' \right] \\ & \leq \frac{7}{2} \int_0^\tau \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [\|\mathbf{z}_{t_n, \tau'}\|^2] d\tau' + \left(\frac{1}{2} \mathbb{E} [\|\hat{\mathbf{y}}_{t_n, \tau}^{(0)}\|^2] + M_s^2 + d \right) \tau \\ & + 3L_s^2 \int_0^\tau \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\left\| \hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(K)}(\omega) - \hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(K-1)}(\omega) \right\|^2 \right] d\tau' + 3 \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [A_{t_n}(\omega) + B_{t_n}(\omega)], \end{aligned}$$

where in the second inequality above, we have used the fact that $s_t^\theta(\cdot)$ is L_s -Lipschitz for any t . By Grönwall's inequality, we have that for any $\tau \in [0, h_n]$

$$\begin{aligned} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [\|\mathbf{z}_{t_n, \tau}\|^2] & \leq e^{\frac{7}{2}\tau} \left[\left(\frac{1}{2} \mathbb{E} [\|\hat{\mathbf{y}}_{t_n, \tau}^{(0)}\|^2] + M_s^2 + d \right) \tau \right] \\ & + 3e^{\frac{7}{2}\tau} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [A_{t_n}(\omega) + B_{t_n}(\omega)] \\ & + 3e^{\frac{7}{2}\tau} L_s^2 \int_0^\tau \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\left\| \hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(K)}(\omega) - \hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(K-1)}(\omega) \right\|^2 \right] d\tau'. \end{aligned} \quad (\text{B.10})$$

By assumption, $\hat{\mathbf{y}}_{t_n, \tau}^{(0)} = \hat{\mathbf{y}}_{t_n}$ follows the distribution of $\tilde{\mathbf{x}}_{t_n} \sim \tilde{p}_{t_n}$, which allows us to bound the second moment of $\hat{\mathbf{y}}_{t_n}$ for any $n \in [0 : N]$ by Lemma A.8:

$$\mathbb{E} [\|\hat{\mathbf{y}}_{t_n}\|^2] = \mathbb{E} [\|\tilde{\mathbf{x}}_{t_n}\|^2] \leq 2d.$$

Substituting (A.5) into (B.10) then yields that for any $\tau \in [0, h_n]$

$$\begin{aligned} \mathbb{E}_{\omega \sim q|\mathcal{F}_{t_n}} [\|\mathbf{z}_{t_n, \tau}\|^2] & \leq \tau e^{\frac{7}{2}\tau} (M_s^2 + 2d) + 3e^{\frac{7}{2}\tau} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [A_{t_n}(\omega) + B_{t_n}(\omega)] \\ & + 3\tau e^{\frac{7}{2}\tau} L_s^2 \sup_{\tau' \in [0, h_n]} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\left\| \hat{\mathbf{y}}_{t_n, \tau'}^{(K)}(\omega) - \hat{\mathbf{y}}_{t_n, \tau'}^{(K-1)}(\omega) \right\|^2 \right]. \end{aligned}$$

Taking supremum with respect to $\tau \in [0, h_n]$ on both sides above completes our proof. \square

As utilized in the proof of the existence of solutions of SDEs, the following lemma demonstrates the exponential convergence of the iteration defined in (B.2).

Lemma B.6 (Exponential convergence of Picard iteration in PIADM-SDE). *For any $n \in [0, N]$, suppose the initialization $\hat{\mathbf{y}}_{t_n}$ in (B.3) of the auxiliary process (B.2) follows the distribution of $\tilde{\mathbf{x}}_{t_n} \sim \tilde{p}_{t_n}$, then the two ending terms $\hat{\mathbf{y}}_{t_n, \tau}^{(K)}$ and $\hat{\mathbf{y}}_{t_n, \tau}^{(K-1)}$ of the sequence $\{\hat{\mathbf{y}}_{t_n, \tau}^{(k)}\}_{k \in [0:K-1]}$ satisfy the following exponential convergence rate*

$$\begin{aligned} & \sup_{\tau \in [0, h_n]} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\left\| \hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) - \hat{\mathbf{y}}_{t_n, \tau}^{(K-1)}(\omega) \right\|_2^2 \right] \\ & \leq \frac{(L_s^2 h_n e^{2h_n})^{K-1} h e^{\frac{7}{2}h_n} (M_s^2 + 2d)}{1 - 3(L_s^2 h_n e^{2h_n})^{K-1} e^{\frac{7}{2}h_n} h_n L_s^2} + \frac{3(L_s^2 h_n e^{2h_n})^{K-1} e^{\frac{7}{2}h_n} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [A_{t_n}(\omega) + B_{t_n}(\omega)]}{1 - 3(L_s^2 h_n e^{2h_n})^{K-1} e^{\frac{7}{2}h_n} h_n L_s^2}. \end{aligned}$$

Proof. For each $\omega \in \Omega$ conditioned on the filtration \mathcal{F}_{t_n} , subtracting (B.2) from the process as defined by

$$d\hat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) = \left[\frac{1}{2} \hat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) + \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(k-1)}(\omega)) \right] d\tau + d\mathbf{w}_{t_n+\tau}(\omega), \quad (\text{B.11})$$

we have

$$\begin{aligned} & d \left(\widehat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) \right) \\ &= \left[\frac{1}{2} \left(\widehat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) \right) + \mathbf{s}_{t_n+g_n(\tau)}^\theta \left(\widehat{\mathbf{y}}_{t_n, g_n(\tau)}^{(k)}(\omega) \right) - \mathbf{s}_{t_n+g_n(\tau)}^\theta \left(\widehat{\mathbf{y}}_{t_n, g_n(\tau)}^{(k-1)}(\omega) \right) \right] d\tau, \end{aligned}$$

where the diffusion term $d\mathbf{w}_{t_n+\tau}$ cancels each other out. Now we may use the formula above to compute derivative $\frac{d}{d\tau'} \left\| \widehat{\mathbf{y}}_{t_n, \tau'}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau'}^{(k)}(\omega) \right\|^2$ explicitly, integrate it from $\tau' = 0$ to τ , and obtain the following inequality

$$\begin{aligned} & \left\| \widehat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) \right\|^2 \\ &= \int_0^\tau 2 \left(\widehat{\mathbf{y}}_{t_n, \tau'}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau'}^{(k)}(\omega) \right)^\top \left(\mathbf{s}_{t_n+g_n(\tau')}^\theta \left(\widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)}(\omega) \right) - \mathbf{s}_{t_n+g_n(\tau')}^\theta \left(\widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k-1)}(\omega) \right) \right) d\tau' \\ &+ \int_0^\tau \left\| \widehat{\mathbf{y}}_{t_n, \tau'}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau'}^{(k)}(\omega) \right\|^2 d\tau' \\ &\leq 2 \int_0^\tau \left\| \widehat{\mathbf{y}}_{t_n, \tau'}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau'}^{(k)}(\omega) \right\|^2 d\tau' \\ &+ \int_0^\tau \left\| \mathbf{s}_{t_n+g_n(\tau')}^\theta \left(\widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)}(\omega) \right) - \mathbf{s}_{t_n+g_n(\tau')}^\theta \left(\widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k-1)}(\omega) \right) \right\|^2 d\tau' \\ &\leq 2 \int_0^\tau \left\| \widehat{\mathbf{y}}_{t_n, \tau'}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau'}^{(k)}(\omega) \right\|^2 d\tau' + L_s^2 \int_0^\tau \left\| \widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)}(\omega) - \widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k-1)}(\omega) \right\|^2 d\tau'. \end{aligned}$$

By Grönwall's inequality, we have

$$\left\| \widehat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) \right\|^2 \leq L_s^2 e^{2\tau} \int_0^\tau \left\| \widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)}(\omega) - \widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k-1)}(\omega) \right\|^2 d\tau'. \quad (\text{B.12})$$

Taking expectation on both sides above further implies that for any $\tau \in [0, h_n]$,

$$\begin{aligned} & \mathbb{E}_{\omega \sim \bar{p} | \mathcal{F}_{t_n}} \left[\left\| \widehat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) \right\|^2 \right] \\ &\leq L_s^2 e^{2\tau} \int_0^\tau \mathbb{E}_{\omega \sim \bar{p} | \mathcal{F}_{t_n}} \left[\left\| \widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)}(\omega) - \widehat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k-1)}(\omega) \right\|^2 \right] d\tau' \\ &\leq L_s^2 \tau e^{2\tau} \sup_{\tau' \in [0, \tau]} \mathbb{E}_{\omega \sim \bar{p} | \mathcal{F}_{t_n}} \left[\left\| \widehat{\mathbf{y}}_{t_n, \tau'}^{(k)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau'}^{(k-1)}(\omega) \right\|^2 \right]. \end{aligned} \quad (\text{B.13})$$

Furthermore, we take supremum over $\tau \in [0, h_n]$ on both sides above and iterate (B.12) over $k \in \mathbb{N}$, which indicates

$$\begin{aligned} & \sup_{\tau \in [0, h_n]} \mathbb{E}_{\omega \sim \bar{p} | \mathcal{F}_{t_n}} \left[\left\| \widehat{\mathbf{y}}_{t_n, \tau}^{(k+1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) \right\|^2 \right] \\ &\leq L_s^2 h_n e^{2h_n} \sup_{\tau \in [0, h_n]} \mathbb{E}_{\omega \sim \bar{p} | \mathcal{F}_{t_n}} \left[\left\| \widehat{\mathbf{y}}_{t_n, \tau}^{(k)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(k-1)}(\omega) \right\|^2 \right] \\ &\leq (L_s^2 h_n e^{2h_n})^k \sup_{\tau \in [0, h_n]} \mathbb{E} \left[\left\| \widehat{\mathbf{y}}_{t_n, \tau}^{(1)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(0)}(\omega) \right\|^2 \right] \\ &\leq (L_s^2 h_n e^{2h_n})^k h e^{\frac{7}{2}h_n} (M_s^2 + 2d) + 3 (L_s^2 h_n e^{2h_n})^k e^{\frac{7}{2}h_n} \mathbb{E}_{\omega \sim \bar{p} | \mathcal{F}_{t_n}} [A_{t_n}(\omega) + B_{t_n}(\omega)] \\ &+ 3 (L_s^2 h_n e^{2h_n})^k e^{\frac{7}{2}h_n} h_n L_s^2 \sup_{\tau \in [0, h_n]} \mathbb{E}_{\omega \sim \bar{p} | \mathcal{F}_{t_n}} \left[\left\| \widehat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) - \widehat{\mathbf{y}}_{t_n, \tau}^{(K-1)}(\omega) \right\|^2 \right], \end{aligned} \quad (\text{B.14})$$

where the last inequality follows from Lemma B.5. By rearranging the inequality above, setting $k = K - 1$ and using the assumption that $L_s^2 h_n e^{2h_n} \ll 1$, we obtain

$$\begin{aligned} & \sup_{\tau \in [0, h_n]} \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\left\| \hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) - \hat{\mathbf{y}}_{t_n, \tau}^{(K-1)}(\omega) \right\|^2 \right] \\ & \leq \frac{(L_s^2 h_n e^{2h_n})^{K-1} h e^{\frac{7}{2} h_n} (M_s^2 + 2d) + 3 (L_s^2 h_n e^{2h_n})^{K-1} e^{\frac{7}{2} h_n} \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} [A_{t_n}(\omega) + B_{t_n}(\omega)]}{1 - 3 (L_s^2 h_n e^{2h_n})^{K-1} e^{\frac{7}{2} h_n} h_n L_s^2}, \end{aligned} \quad (\text{B.15})$$

as desired. \square

The following lemma from [107] bounds the expectation of the term $A_{t_n}(\omega)$ in (B.7):

Lemma B.7 ([107, Section 3.1]). *We have*

$$\mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} [A_{t_n}(\omega)] \lesssim \epsilon d h_n, \quad \text{for } n \in [0 : N - 2], \text{ and } \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} [A_{t_{N-1}}(\omega)] \lesssim \epsilon d \log \eta^{-1},$$

where η is the parameter for early stopping.

Proof. Notice that

$$\begin{aligned} & \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} [A_{t_n}(\omega)] \\ & = \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\int_0^{h_n} \left\| \nabla \log \tilde{p}_{t_n + g_n(\tau)}(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) - \nabla \log \tilde{p}_{t_n + \tau}(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) \right\|^2 d\tau \right] \\ & = \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\sum_{m=0}^{M_n} \int_{\tau_{n,m}}^{\tau_{n,m+1}} \left\| \nabla \log \tilde{p}_{t_n + \tau_{n,m}}(\hat{\mathbf{y}}_{t_n, \tau_{n,m}}^{(K)}(\omega)) - \nabla \log \tilde{p}_{t_n + \tau}(\hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega)) \right\|^2 d\tau \right], \\ & = \sum_{m=0}^{M_n} \int_{\tau_{n,m}}^{\tau_{n,m+1}} \mathbb{E}_{\omega \sim \tilde{p} | \mathcal{F}_{t_n}} \left[\left\| \nabla \log \tilde{p}_{t_n + \tau_{n,m}}(\tilde{\mathbf{x}}_{t_n + \tau}(\omega)) - \nabla \log \tilde{p}_{t_n + \tau}(\tilde{\mathbf{x}}_{t_n + \tau}(\omega)) \right\|^2 \right] d\tau, \end{aligned}$$

where for the last equality, we use the fact that the process $\hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega)$ follows the backward SDE with the true score function under the measure \tilde{p} . In the following, we drop the superscript $\omega \sim \tilde{p} | \mathcal{F}_{t_n}$ of the expectation for simplicity.

By Lemma A.6 and A.7, we have

$$\begin{aligned} & \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_n + \tau_{n,m}}(\tilde{\mathbf{x}}_{t_n + \tau}(\omega)) - \nabla \log \tilde{p}_{t_n + \tau}(\tilde{\mathbf{x}}_{t_n + \tau}(\omega)) \right\|^2 \right] \\ & \leq \int_0^\tau \left(\frac{1}{2} \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_n + \tau_{n,m}}(\tilde{\mathbf{x}}_{t_n + \tau_{n,m}}(\omega)) \right\|^2 \right] + \mathbb{E} \left[\left\| \nabla^2 \log \tilde{p}_{t_n + \tau'}(\tilde{\mathbf{x}}_{t_n + \tau'}(\omega)) \right\|_F^2 \right] \right) d\tau' \\ & \leq \int_0^\tau \left(\frac{1}{2} d \bar{\sigma}_{\tau'}^{-2} + d \bar{\sigma}_{\tau'}^{-4} \right) d\tau' + \left(\bar{\sigma}_{t_n + \tau_{n,m}}^{-4} \mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n + \tau_{n,m}} \right] - \bar{\sigma}_{t_n + \tau}^{-4} \mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n + \tau} \right] \right), \end{aligned}$$

Now noticing that

$$\bar{\sigma}_t^2 = \sigma_{T-t}^2 \lesssim T - t,$$

we further have

$$\begin{aligned} & \int_{\tau_{n,m}}^{\tau_{n,m+1}} \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_n + \tau_{n,m}}(\tilde{\mathbf{x}}_{t_n + \tau}(\omega)) - \nabla \log \tilde{p}_{t_n + \tau}(\tilde{\mathbf{x}}_{t_n + \tau}(\omega)) \right\|^2 \right] d\tau \\ & \lesssim \int_{\tau_{n,m}}^{\tau_{n,m+1}} \int_0^{\tau'} \frac{d}{(T - t_n - \tau_{n,m+1})^2} d\tau' d\tau + \frac{\epsilon_{n,m} \left(\mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n + \tau_{n,m}} \right] - \mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n + \tau_{n,m+1}} \right] \right)}{(T - t_n - \tau_{n,m})^2} \\ & \lesssim d \frac{\epsilon_{n,m}^2}{(T - t_n - \tau_{n,m+1})^2} + \frac{\epsilon \left(\mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n + \tau_{n,m}} \right] - \mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n + \tau_{n,m+1}} \right] \right)}{T - t_n - \tau_{n,m}}, \end{aligned}$$

and thus

$$\begin{aligned}
& \sum_{m=0}^{M_n} \int_{\tau_{n,m}}^{\tau_{n,m+1}} \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_n+\tau_{n,m}}(\tilde{\mathbf{x}}_{t_n+\tau}(\omega)) - \nabla \log \tilde{p}_{t_n+\tau}(\tilde{\mathbf{x}}_{t_n+\tau}(\omega)) \right\|^2 \right] d\tau \\
& \lesssim d \sum_{m=0}^{M_n} \frac{\epsilon_{n,m}^2}{(T-t_n-\tau_{n,m+1})^2} + \sum_{m=0}^{M_n} \frac{\epsilon}{T-t_n-\tau_{n,m}} \left(\mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n+\tau_{n,m}} \right] - \mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n+\tau_{n,m+1}} \right] \right) \\
& \leq d\epsilon^2 M_n + \frac{\epsilon \mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n+\tau_{n,0}} \right]}{T-t_n-\tau_{n,0}} + \sum_{m=0}^{M_n} \frac{\epsilon \epsilon_{n,m} \mathbb{E} \left[\text{tr} \tilde{\Sigma}_{t_n+\tau_{n,m}} \right]}{(T-t_n-\tau_{n,m+1})(T-t_n-\tau_{n,m})} \\
& \leq d\epsilon^2 M_n + \epsilon d + d\epsilon^2 M_n \lesssim d\epsilon^2 M_n.
\end{aligned}$$

For $n \in [0, N-2]$, we have $M_n \epsilon = h_n$ and thus $\mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} [A_{t_n}(\omega)] \lesssim \epsilon d h_n$, and for $n = N-1$, we have

$$M_N \lesssim \int_{\eta}^h \frac{1}{\epsilon \tau} d\tau = \log \eta^{-1} \epsilon^{-1}$$

and thus $\mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} [A_{t_{N-1}}(\omega)] \lesssim \epsilon^2 d M_n \lesssim \epsilon d \log \eta^{-1}$. \square

B.3 Overall Error Bound

Proof of Theorem 3.3. We first continue the computation in (B.6) and (B.7):

$$\begin{aligned}
& D_{\text{KL}}(\tilde{p}_{t_{n+1}} \parallel \hat{q}_{t_{n+1}}) \leq D_{\text{KL}}(\tilde{p}_{t_n} \parallel \hat{q}_{t_n}) + \mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} \left[\frac{1}{2} \int_0^{h_n} \|\delta_{t_n}(\tau, \omega)\|^2 d\tau \right] \\
& \leq D_{\text{KL}}(\tilde{p}_{t_n} \parallel \hat{q}_{t_n}) + 3\mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} [A_{t_n}(\omega) + B_{t_n}(\omega)] \\
& + 3\mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} \left[\int_0^{h_n} \left\| \mathbf{s}_{t_n+g_n(\tau)}^{\theta}(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega)) - \mathbf{s}_{t_n+g_n(\tau)}^{\theta}(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K-1)}(\omega)) \right\|^2 d\tau \right] \\
& \leq D_{\text{KL}}(\tilde{p}_{t_n} \parallel \hat{q}_{t_n}) + 3\mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} \left[A_{t_n}(\omega) + B_{t_n}(\omega) + L_s^2 \int_0^{h_n} \left\| \hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K)}(\omega) - \hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(K-1)}(\omega) \right\|^2 d\tau \right] \\
& \leq D_{\text{KL}}(\tilde{p}_{t_n} \parallel \hat{q}_{t_n}) + 3\mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} \left[A_{t_n}(\omega) + B_{t_n}(\omega) + h_n L_s^2 \sup_{\tau \in [0, h_n]} \left\| \hat{\mathbf{y}}_{t_n, \tau}^{(K)}(\omega) - \hat{\mathbf{y}}_{t_n, \tau}^{(K-1)}(\omega) \right\|^2 \right].
\end{aligned}$$

Then plugging in the result of Lemma B.6, we have

$$\begin{aligned}
& D_{\text{KL}}(\tilde{p}_{t_{n+1}} \parallel \hat{q}_{t_{n+1}}) \\
& \leq D_{\text{KL}}(\tilde{p}_{t_n} \parallel \hat{q}_{t_n}) + 3\mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} [A_{t_n}(\omega) + B_{t_n}(\omega)] + 3h_n L_s^2 \frac{(L_s^2 h_n e^{2h_n})^{K-1} h e^{\frac{7}{2}h_n} (M_s^2 + 2d)}{1 - 3(L_s^2 h_n e^{2h_n})^{K-1} e^{\frac{7}{2}h_n} h_n L_s^2} \\
& + h_n L_s^2 \frac{9(L_s^2 h_n e^{2h_n})^{K-1} e^{\frac{7}{2}h_n} \mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} [A_{t_n}(\omega) + B_{t_n}(\omega)]}{1 - 3(L_s^2 h_n e^{2h_n})^{K-1} e^{\frac{7}{2}h_n} h_n L_s^2} \\
& \lesssim D_{\text{KL}}(\tilde{p}_{t_n} \parallel \hat{q}_{t_n}) + \frac{1 + e^{-K} h_n e^{h_n}}{1 - e^{-K} h_n e^{h_n}} \mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} [A_{t_n}(\omega) + B_{t_n}(\omega)] + e^{-K} h_n^2 e^{h_n} d \\
& \lesssim D_{\text{KL}}(\tilde{p}_{t_n} \parallel \hat{q}_{t_n}) + \mathbb{E}_{\omega \sim \tilde{p}|_{\mathcal{F}_{t_n}}} [A_{t_n}(\omega) + B_{t_n}(\omega)] + e^{-K} h_n^2 e^{h_n} d,
\end{aligned}$$

where we used the assumption that $L_s^2 h_n e^{\frac{7}{2}h_n} \ll 1$.

The term $\sum_{n=0}^{N-1} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [B_{t_n}(\omega)]$ is bounded by Assumption 3.1 as

$$\begin{aligned} & \sum_{n=0}^{N-1} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [B_{t_n}(\omega)] \\ & \leq \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\sum_{n=0}^{N-1} \int_0^{h_n} \left\| \mathbf{s}_{t_n+g_n(\tau)}^\theta(\hat{\mathbf{y}}_{t_n,g_n(\tau)}^{(K)}(\omega)) - \nabla \log \tilde{p}_{t_n+g_n(\tau)}(\hat{\mathbf{y}}_{t_n,g_n(\tau)}^{(K)}(\omega)) \right\|^2 d\tau \right] \\ & = \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\sum_{n=0}^{N-1} \sum_{m=0}^{M_n-1} \epsilon_{n,m} \left\| \mathbf{s}_{t_n+\tau_{n,m}}^\theta(\hat{\mathbf{y}}_{t_n,\tau_{n,m}}^{(K)}(\omega)) - \nabla \log \tilde{p}_{t_n+\tau_{n,m}}(\hat{\mathbf{y}}_{t_n,\tau_{n,m}}^{(K)}(\omega)) \right\|^2 \right] \\ & = \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} \left[\sum_{n=0}^{N-1} \sum_{m=0}^{M_n-1} \epsilon_{n,m} \left\| \mathbf{s}_{t_n+\tau_{n,m}}^\theta(\tilde{\mathbf{x}}_{t_n+\tau}(\omega)) - \nabla \log \tilde{p}_{t_n+\tau_{n,m}}(\tilde{\mathbf{x}}_{t_n+\tau}(\omega)) \right\|^2 \right] \leq \delta_2^2, \end{aligned}$$

where the last equality is because the process $\hat{\mathbf{y}}_{t_n,\tau}^{(K)}(\omega)$ under measure \tilde{p} follows the backward SDE (B.5).

Thus, by Theorem A.1 and plugging in the iteration relations above

$$\begin{aligned} D_{\text{KL}}(p_\eta \| \hat{q}_{t_N}) &= D_{\text{KL}}(\tilde{p}_{t_N} \| \hat{q}_{t_N}) \\ &\leq D_{\text{KL}}(\tilde{p}_0 \| \hat{q}_0) + \sum_{n=0}^{N-1} \left(\mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [A_{t_n}(\omega) + B_{t_n}(\omega)] + e^{-K} h_n^2 e^{h_n} d \right) \\ &\leq D_{\text{KL}}(\tilde{p}_0 \| \hat{q}_0) + \sum_{n=0}^{N-2} \epsilon d h_n + \epsilon d \log \eta^{-1} + \sum_{n=0}^{N-1} \mathbb{E}_{\omega \sim \tilde{p}|\mathcal{F}_{t_n}} [B_{t_n}(\omega)] + e^{-K} h_n^2 e^{h_n} d N \\ &\leq d e^{-T} + \epsilon d (T + \log \eta^{-1}) + \delta_2^2 + e^{-K} d T \leq d e^{-T} + \epsilon d T + \delta^2 + e^{-K} d T, \end{aligned}$$

as $T \gtrsim \log \eta^{-1}$, $h_n \lesssim 1$, and $\delta_2 \lesssim \delta$, and then it is straightforward to see that the following choices of parameters

$$\begin{aligned} T &= \mathcal{O}(\log(d\delta^{-2})), \quad h = \Theta(1), \quad N = \mathcal{O}(\log(d\delta^{-2})), \\ \epsilon &= \Theta(d^{-1}\delta^2 \log^{-1}(d\delta^{-2})), \quad M = \mathcal{O}(d\delta^{-2} \log(d\delta^{-2})), \\ K &= \tilde{\mathcal{O}}(\log(d\delta^{-2})), \end{aligned}$$

would yield an overall error of $\mathcal{O}(\delta^2)$. □

C Details of Probability Flow ODE Implementation

In this section, we provide the details of the parallelized algorithm for the probability flow ODE formulation of diffusion models. We first introduce the algorithm and define the necessary notations, then discuss the error analysis during the predictor and corrector steps, respectively, and finally provide the proof of Theorem 3.5.

C.1 Algorithm

In the parallelized inference algorithm for diffusion models in the probability flow ODE formulation, we adopt the same discretization scheme as in Section 3.1.1 and the exponential integrator for all updating rules. For each block, we first run a *predictor step*, which consists of running the probability flow ODE in parallel. Then we run a *corrector step*, which runs an underdamped Langevin dynamics in parallel to correct the distribution of the samples. The algorithm is summarized In Algorithm 2.

Parallelized Predictor Step The parallelization strategies in the predictor step are similar to those in the SDE algorithm (Algorithm 1). The only difference here is that instead of applying Picard iteration to the backward SDE as in (3.2), we apply Picard iteration to the probability flow ODE as in (C.3), which does not require i.i.d. samples from standard Gaussian distribution. As shown in Lemma C.3, the update rule in the predictor step (C.1) in Algorithm 2 is equivalent to running

Algorithm 2: PIADM-ODE

Input: $\hat{\mathbf{y}}_0 \sim \hat{q}_0 = \mathcal{N}(0, \mathbf{I}_d)$, a discretization scheme $(T, (h_n)_{n=1}^N$ and $(\tau_{n,m})_{n \in [1:N], m \in [0:M_n]}$) satisfying (3.1), parameters for the corrector step $(T^\dagger, N^\dagger, h^\dagger, M^\dagger, \epsilon^\dagger)$, the depth of iteration K and K^\dagger , the learned NN-based score $\mathbf{s}_t^\theta(\cdot)$.

Output: A sample $\hat{\mathbf{y}}_T \sim \hat{q}_T \approx \tilde{p}_T$.

```
1 for  $n = 0$  to  $N - 1$  do
2   ▷ Predictor Step (Section C.2)
3    $\hat{\mathbf{y}}_{t_n, \tau_{n,m}}^{(0)} \leftarrow \hat{\mathbf{y}}_{t_n}$  for  $m \in [0 : M_n]$ ;
4   for  $k = 1$  to  $K$  do
5      $\hat{\mathbf{y}}_{t_n, 0}^{(k)} \leftarrow \hat{\mathbf{y}}_{t_n}$ ;
6     for  $m = 1$  to  $M_n$  in parallel do
7       
$$\begin{aligned} \hat{\mathbf{y}}_{t_n, \tau_{n,m}}^{(k)} \leftarrow & \frac{1}{2} e^{\frac{\tau_{n,m}}{2}} \hat{\mathbf{y}}_{t_n, 0}^{(k-1)} \\ & + \frac{1}{2} \sum_{j=0}^{m-1} e^{\frac{\tau_{n,m} - \tau_{n,j+1}}{2}} (e^{\epsilon_{n,j}} - 1) \mathbf{s}_{t_n + \tau_{n,j}}^\theta(\hat{\mathbf{y}}_{t_n, \tau_{n,j}}^{(k-1)}) \end{aligned} \quad (\text{C.1})$$

8     end
9   end
10  ▷ Corrector Step (Section C.3)
11   $\hat{\mathbf{u}}_{t_n, 0}^{(0)} \leftarrow \hat{\mathbf{y}}_{t_n, h_n}^{(K)}$  and  $\hat{\mathbf{v}}_{t_n, 0}^{(0)} \sim \mathcal{N}(0, \mathbf{I}_d)$ ;
12  for  $n^\dagger = 0$  to  $N^\dagger - 1$  do
13     $(\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, m^\dagger \epsilon^\dagger}^{(0)}, \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, m^\dagger \epsilon^\dagger}^{(0)}) \leftarrow (\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger}, \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger})$  for  $m^\dagger \in [0 : M^\dagger]$ ;
14    for  $k^\dagger = 1$  to  $K^\dagger$  do
15       $(\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, 0}^{(k^\dagger)}, \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, 0}^{(k^\dagger)}) \leftarrow (\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger}, \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger})$ ;
16       $\boldsymbol{\xi}_{j^\dagger} \sim \mathcal{N}\left(\mathbf{0}, 2\gamma(1 + \gamma^{-2})(1 - e^{-\gamma \epsilon^\dagger})^2 e^{-2\gamma((M^\dagger - j^\dagger + 1)\epsilon^\dagger)} \mathbf{I}_d\right)$  for
17         $j^\dagger \in [0 : M^\dagger]$ ;
18      for  $m^\dagger = 1$  to  $M^\dagger$  in parallel do
19        
$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, m^\dagger \epsilon^\dagger}^{(k^\dagger)} \\ \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, m^\dagger \epsilon^\dagger}^{(k^\dagger)} \end{bmatrix} \leftarrow & \mathbf{G}(m^\dagger \epsilon^\dagger) \begin{bmatrix} \hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, 0}^{(k^\dagger-1)} \\ \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, 0}^{(k^\dagger-1)} \end{bmatrix} \\ & + \sum_{j^\dagger=0}^{m^\dagger-1} \mathbf{G}((m^\dagger - j^\dagger - 1)\epsilon^\dagger) (\mathbf{I}_d - \mathbf{G}(\epsilon^\dagger)) \begin{bmatrix} \mathbf{0} \\ \mathbf{s}_{t_n+1}^\theta(\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, j^\dagger \epsilon^\dagger}^{(k^\dagger-1)}) \end{bmatrix} \\ & + \sum_{j^\dagger=0}^{m^\dagger-1} \mathbf{G}((m^\dagger - j^\dagger - 1)\epsilon^\dagger) \begin{bmatrix} \mathbf{0} \\ \boldsymbol{\xi}_{j^\dagger} \end{bmatrix}; \end{aligned} \quad (\text{C.2})$$

19      end
20    end
21     $(\hat{\mathbf{u}}_{t_n, (n^\dagger+1)h^\dagger}, \hat{\mathbf{v}}_{t_n, (n^\dagger+1)h^\dagger}) \leftarrow (\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger}^{(K^\dagger)}, \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger}^{(K^\dagger)})$ ;
22  end
23   $\hat{\mathbf{y}}_{t_{n+1}} \leftarrow \hat{\mathbf{u}}_{t_n, T^\dagger}$ ;
24 end
```

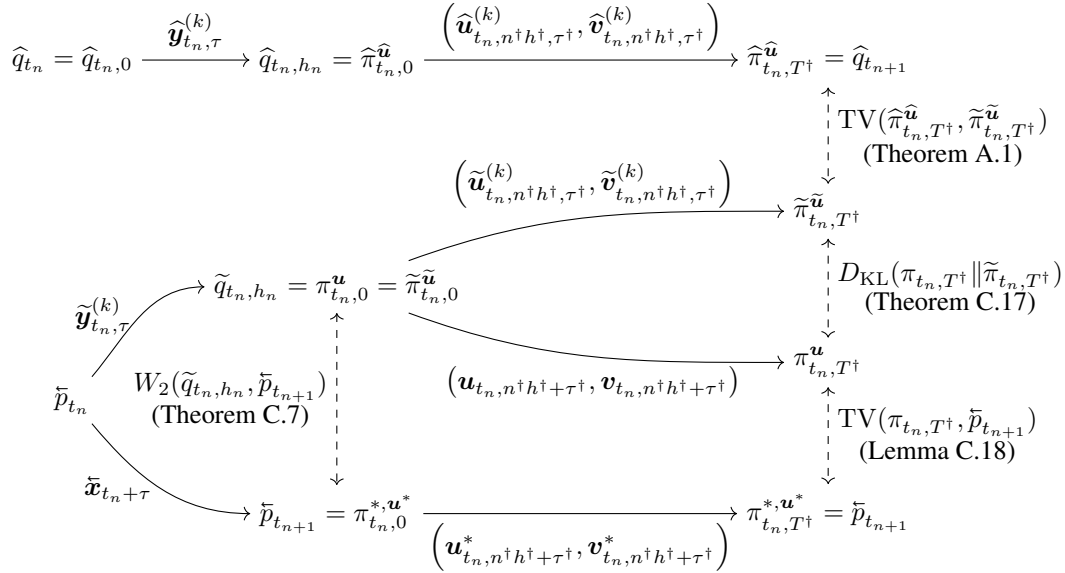


Figure 2: Illustration of the proof pipeline of Theorem 3.5 for PIADM-ODE within the n -th block.

the auxiliary predictor process (C.3). The auxiliary predictor process takes in the result from the previous corrector step (or the initialization if $n = 0$) and outputs $\hat{\mathbf{y}}_{t_n,h_n}^{(K)}$ as the initialization for the next corrector step.

Parallelized Corrector Step The parallelization of the underdamped Langevin dynamics is similar to that mentioned in Section 2.2. Given a sample resulting from the predictor step, we initialize the auxiliary corrector process (Definition C.8) which is an underdamped Langevin dynamics with the initialization $\hat{\mathbf{u}}_{t_n,0} = \mathbf{y}_{t_n,h_n}^{(K)}$ and the augmented variable $\hat{\mathbf{v}}_{t_n,0} \sim \mathcal{N}(0, \mathbf{I}_d)$ representing the momentum.

We run the underdamped Langevin dynamics for time T^\dagger , which is set to be of order $\Omega(1)$ so that it is large enough to correct the distribution of the samples (*cf.* Lemma C.18) while being comparably short to ensure numerical stability (*cf.* Theorem C.17). Following a similar strategy as in Section 2.2 and in Algorithm 1, we further divide the time horizon T^\dagger into N^\dagger blocks with step size h^\dagger , and for each block the block length h^\dagger into M^\dagger steps with step size ϵ^\dagger . Within each block, we run the underdamped Langevin dynamics in parallel for K^\dagger iterations. As shown in Lemma C.9, the update rule in the corrector step (C.2) in Algorithm 2 is equivalent to running the auxiliary corrector process (C.11).

In the following subsections, we proceed to provide theoretical guarantees for the algorithm.

C.2 Parallelized Predictor Step

Definition C.1 (Auxiliary Predictor Process). *For any $n \in [0 : N - 1]$, we define the auxiliary predictor process $(\hat{\mathbf{y}}_{t_n,\tau}^{(k)})_{\tau \in [0, h_n]}$ as the solution to the following ODE recursively for $k \in [0 : K - 1]$:*

$$d\hat{\mathbf{y}}_{t_n,\tau}^{(k+1)} = \left[\frac{1}{2} \hat{\mathbf{y}}_{t_n,\tau}^{(k+1)} + \frac{1}{2} \mathbf{s}_{t_n+g_n(\tau)}^\theta \left(\hat{\mathbf{y}}_{t_n,g_n(\tau)}^{(k)} \right) \right] d\tau, \quad (\text{C.3})$$

with the initial condition

$$\hat{\mathbf{y}}_{t_n,\tau}^{(0)} \equiv \hat{\mathbf{y}}_{t_n} \text{ for } \tau \in [0, h_n], \quad \text{and} \quad \mathbf{y}_{t_n,0}^{(k)} \equiv \hat{\mathbf{y}}_{t_n} \text{ for } k \in [1 : K] \quad (\text{C.4})$$

where $\hat{\mathbf{y}}_{t_n} = \hat{\mathbf{u}}_{t_{n-1},N^\dagger h^\dagger}$ if $n \in [1 : N - 1]$ and $\hat{\mathbf{y}}_{t_0} \sim \mathcal{N}(0, \mathbf{I}_d)$. We will also denote the probability distribution of $\hat{\mathbf{y}}_{t_n,\tau}^{(K)}$ as $\hat{q}_{t_n,\tau}$.

Definition C.2 (Interpolating Process). For any $n \in [0 : N - 1]$, we define the interpolating process $(\tilde{\mathbf{y}}_{t_n, \tau}^{(k)})_{\tau \in [0, h_n]}$ as the solution to the following ODE recursively for $k \in [0 : K - 1]$:

$$d\tilde{\mathbf{y}}_{t_n, \tau}^{(k+1)} = \left[\frac{1}{2} \tilde{\mathbf{y}}_{t_n, \tau}^{(k+1)} + \frac{1}{2} \mathbf{s}_{t_n + g_n(\tau)}^\theta \left(\tilde{\mathbf{y}}_{t_n, g_n(\tau)}^{(k)} \right) \right] d\tau, \quad (\text{C.5})$$

with initial condition

$$\tilde{\mathbf{y}}_{t_n, \tau}^{(0)} \equiv \tilde{\mathbf{y}}_{t_n, 0}^{(0)} \quad \text{for } \tau \in [0, h_n], \text{ and } \tilde{\mathbf{y}}_{t_n, 0}^{(k)} \equiv \tilde{\mathbf{y}}_{t_n, 0}^{(0)} \quad \text{for } k \in [1 : K],$$

where $\tilde{\mathbf{y}}_{t_n, 0}^{(0)} \sim \tilde{p}_{t_n}$. We will also denote the probability distribution of $\tilde{\mathbf{y}}_{t_n, \tau}^{(K)}$ as $\tilde{q}_{t_n, \tau}$.

Similar to the equivalence between (3.4) and (B.2), we have the following lemma:

Lemma C.3 (Equivalence between (C.1) and (C.3)). For any $n \in [0 : N - 1]$, the update rule (C.1) in Algorithm 2 is equivalent to the exact solution of (C.3) for any $k \in [0 : K - 1]$ and $\tau \in [0, h_n]$.

Proof. Rewriting (C.3) and multiplying $e^{-\frac{\tau}{2}}$ on both sides yield

$$d \left[e^{-\frac{\tau}{2}} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} \right] = e^{-\frac{\tau}{2}} \left[d\hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \frac{1}{2} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} d\tau \right] = \frac{e^{-\frac{\tau}{2}}}{2} \mathbf{s}_{t_n + g_n(\tau)}^\theta \left(\hat{\mathbf{y}}_{t_n, g_n(\tau)}^{(k)} \right) d\tau$$

Integrating on both sides from 0 to τ implies

$$\begin{aligned} e^{-\frac{\tau}{2}} \hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \hat{\mathbf{y}}_{t_n, 0}^{(k+1)} &= \int_0^\tau \frac{e^{-\frac{\tau'}{2}}}{2} \mathbf{s}_{t_n + g_n(\tau')}^\theta \left(\hat{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)} \right) d\tau' \\ &= \frac{1}{2} \sum_{m=0}^{M_n} \int_{\tau \wedge t_{n, m}}^{\tau \wedge \tau_{n, m+1}} e^{-\frac{\tau'}{2}} \mathbf{s}_{t_n + \tau_{n, m}}^\theta \left(\hat{\mathbf{y}}_{t_n, \tau_{n, m}}^{(k)} \right) d\tau' \\ &= \sum_{m=0}^{M_n} \left(e^{-\frac{\tau \wedge \tau_{n, m}}{2}} - e^{-\frac{\tau \wedge \tau_{n, m+1}}{2}} \right) \mathbf{s}_{t_n + \tau_{n, m}}^\theta \left(\hat{\mathbf{y}}_{t_n, \tau_{n, m}}^{(k)} \right), \end{aligned}$$

and then multiplying $e^{\frac{\tau}{2}}$ on both sides above yields

$$\hat{\mathbf{y}}_{t_n, \tau}^{(k+1)} = e^{\frac{\tau}{2}} \hat{\mathbf{y}}_{t_n, 0}^{(k+1)} + \sum_{m=0}^{M_n} \left(e^{\frac{\tau \wedge \tau_{n, m+1} - \tau \wedge \tau_{n, m}}{2}} - 1 \right) e^{\frac{0 \vee (\tau - \tau_{n, m+1})}{2}} \mathbf{s}_{t_n + \tau_{n, m}}^\theta \left(\hat{\mathbf{y}}_{t_n, \tau_{n, m}}^{(k)} \right).$$

Plugging in $\tau = \tau_{n, m}$ gives us (C.1), as desired. \square

Lemma C.4 (Error between the interpolating process and the true process). Under the Picard iteration, we have that the ending process $\{\hat{\mathbf{y}}_{t_n, \tau}^{(K)}\}_{\tau \in [0, h_n]}$ satisfies the following exponential convergence rate

$$\sup_{\tau \in [0, h_n]} \mathbb{E} \left[\left\| \hat{\mathbf{y}}_{t_n, \tau}^{(K)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] \leq 3d \left(\frac{h_n^2 e^{h_n + \frac{3}{2}} L_s^2}{2} \right)^K + \frac{e^{h_n + \frac{3}{2}} h_n / 2}{1 - h_n^2 e^{h_n + \frac{3}{2}} L_s^2 / 2} \left(h_n \delta_\infty^2 + \mathbb{E}[D_{t_n}] \right),$$

where

$$D_{t_n} := \int_0^{h_n} \left\| \mathbf{s}_{t_n + g_n(\tau')}^\theta \left(\tilde{\mathbf{x}}_{t_n + g_n(\tau')} \right) - \mathbf{s}_{t_n + \tau'}^\theta \left(\tilde{\mathbf{x}}_{t_n + \tau'} \right) \right\|^2 d\tau'.$$

Proof. Recall that the backward true process $\{\tilde{\mathbf{x}}_{t_n + \tau}\}_{\tau \in [0, h_n]}$ satisfies the following backward SDE within one block

$$d\tilde{\mathbf{x}}_{t_n + \tau} = \left[\frac{1}{2} \tilde{\mathbf{x}}_{t_n + \tau} + \frac{1}{2} \nabla \log \tilde{p}_{t_n + \tau}(\tilde{\mathbf{x}}_{t_n + \tau}) \right] d\tau. \quad (\text{C.6})$$

By subtracting (C.6) from (C.5), we obtain that

$$\begin{aligned} \frac{d}{d\tau} \left(\tilde{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau} \right) &= \frac{1}{2} \left[\tilde{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau} \right] \\ &\quad + \frac{1}{2} \left[\mathbf{s}_{t_n + g_n(\tau)}^\theta \left(\tilde{\mathbf{y}}_{t_n, g_n(\tau)}^{(k)} \right) - \mathbf{s}_{t_n + g_n(\tau)}^\theta \left(\tilde{\mathbf{x}}_{t_n + g_n(\tau)} \right) \right] \\ &\quad + \frac{1}{2} \left[\mathbf{s}_{t_n + g_n(\tau)}^\theta \left(\tilde{\mathbf{x}}_{t_n + g_n(\tau)} \right) - \nabla \log \tilde{p}_{t_n + g_n(\tau)}(\tilde{\mathbf{x}}_{t_n + g_n(\tau)}) \right] \\ &\quad + \frac{1}{2} \left[\nabla \log \tilde{p}_{t_n + g_n(\tau)}(\tilde{\mathbf{x}}_{t_n + g_n(\tau)}) - \nabla \log \tilde{p}_{t_n + \tau}(\tilde{\mathbf{x}}_{t_n + \tau}) \right]. \end{aligned} \quad (\text{C.7})$$

Then by

$$\mathrm{d} \left\| \tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right\|^2 = 2 \left(\tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right)^\top \mathrm{d} \left(\tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right),$$

and integrating for $\tau' \in [0, h_n]$, we have

$$\begin{aligned} & \left\| \tilde{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \\ &= \int_0^\tau \left(\tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right)^\top \left(\mathbf{s}_{t_n + g_n(\tau')}^\theta(\tilde{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)}) - \mathbf{s}_{t_n + g_n(\tau')}^\theta(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) \right) \mathrm{d}\tau' \\ &+ \int_0^\tau \left(\tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right)^\top \left(\mathbf{s}_{t_n + g_n(\tau')}^\theta(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) - \nabla \log \tilde{p}_{t_n + g_n(\tau')}(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) \right) \mathrm{d}\tau' \\ &+ \int_0^\tau \left(\tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right)^\top \left(\nabla \log \tilde{p}_{t_n + g_n(\tau')}(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) - \nabla \log \tilde{p}_{t_n + \tau'}(\tilde{\mathbf{x}}_{t_n + \tau'}) \right) \mathrm{d}\tau' \\ &+ \int_0^\tau \left\| \tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right\|^2 \mathrm{d}\tau'. \end{aligned}$$

Using AM-GM inequality and taking expectations on both sides, we further upper bound the summation above as

$$\begin{aligned} & \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] \\ & \leq \left(1 + \frac{3}{2h_n} \right) \int_0^\tau \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right\|^2 \right] \mathrm{d}\tau' \\ & + \frac{h_n}{2} \int_0^\tau \mathbb{E} \left[\left\| \mathbf{s}_{t_n + g_n(\tau')}^\theta(\tilde{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)}) - \mathbf{s}_{t_n + g_n(\tau')}^\theta(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) \right\|^2 \right] \mathrm{d}\tau' \\ & + \frac{h_n}{2} \int_0^\tau \mathbb{E} \left[\left\| \mathbf{s}_{t_n + g_n(\tau')}^\theta(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) - \nabla \log \tilde{p}_{t_n + g_n(\tau')}(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) \right\|^2 \right] \mathrm{d}\tau' \\ & + \frac{h_n}{2} \underbrace{\mathbb{E} \left[\int_0^\tau \left\| \nabla \log \tilde{p}_{t_n + g_n(\tau')}(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) - \nabla \log \tilde{p}_{t_n + \tau'}(\tilde{\mathbf{x}}_{t_n + \tau'}) \right\|^2 \mathrm{d}\tau' \right]}_{\leq D_{t_n}} \\ & \leq \left(1 + \frac{3}{2h_n} \right) \int_0^\tau \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau'}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right\|^2 \right] \mathrm{d}\tau' + \frac{h_n}{2} (\tau \delta_\infty^2 + \mathbb{E}[D_{t_n}]) \\ & + \frac{L_s^2 h_n}{2} \int_0^\tau \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)} - \tilde{\mathbf{x}}_{t_n + g_n(\tau')} \right\|^2 \right] \mathrm{d}\tau', \end{aligned}$$

where the last equality is by Assumption 3.1'.

Applying Grönwall's inequality, we have

$$\begin{aligned} & \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] \\ & \leq \frac{e^{(1 + \frac{3}{2h_n})\tau} L_s^2 h_n}{2} \int_0^\tau \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, g_n(\tau')}^{(k)} - \tilde{\mathbf{x}}_{t_n + g_n(\tau')} \right\|^2 \right] \mathrm{d}\tau' + \frac{e^{(1 + \frac{3}{2h_n})\tau} h_n}{2} (\tau \delta_\infty^2 + \mathbb{E}[D_{t_n}]) \\ & \leq \frac{\tau e^{(1 + \frac{3}{2h_n})\tau} L_s^2 h_n}{2} \sup_{\tau' \in [0, \tau]} \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau'}^{(k)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right\|^2 \right] + \frac{e^{(1 + \frac{3}{2h_n})\tau} h_n}{2} (\tau \delta_\infty^2 + \mathbb{E}[D_{t_n}]), \end{aligned} \tag{C.8}$$

and by taking supremum

$$\begin{aligned} & \sup_{\tau \in [0, h_n]} \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau}^{(k+1)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] \\ & \leq \frac{h_n^2 e^{h_n + \frac{3}{2}} L_s^2}{2} \sup_{\tau' \in [0, \tau]} \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau'}^{(k)} - \tilde{\mathbf{x}}_{t_n + \tau'} \right\|^2 \right] + \frac{e^{h_n + \frac{3}{2}} h_n}{2} (h_n \delta_\infty^2 + \mathbb{E}[D_{t_n}]) \end{aligned} \tag{C.9}$$

Given that constant h_n is sufficiently small, which ensures $L_s^2 h_n e^{\frac{3}{2}h_n} \ll 1$, iterating the above inequality for $k \in [0 : K - 1]$ gives us that

$$\begin{aligned} & \sup_{\tau \in [0, h_n]} \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau}^{(K)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] \\ & \leq \left(\frac{h_n^2 e^{h_n + \frac{3}{2}} L_s^2}{2} \right)^K \sup_{\tau \in [0, h_n]} \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau}^{(0)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] + \frac{e^{h_n + \frac{3}{2}} h_n / 2}{1 - h_n^2 e^{h_n + \frac{3}{2}} L_s^2 / 2} (h_n \delta_\infty^2 + \mathbb{E}[D_{t_n}]), \end{aligned}$$

Notice that by Lemma A.8, we have

$$\mathbb{E} \left[\left\| \hat{\mathbf{y}}_{t_n, \tau}^{(0)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] = \mathbb{E} \left[\left\| \tilde{\mathbf{x}}_{t_n} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] \leq 3d,$$

substituting which into (C.9) then gives us that

$$\sup_{\tau \in [0, h_n]} \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau}^{(K)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] \leq 3d \left(\frac{h_n^2 e^{h_n + \frac{3}{2}} L_s^2}{2} \right)^K + \frac{e^{h_n + \frac{3}{2}} h_n / 2}{1 - h_n^2 e^{h_n + \frac{3}{2}} L_s^2 / 2} (h_n \delta_\infty^2 + \mathbb{E}[D_{t_n}]),$$

as desired. \square

Now it remains to bound C_{t_n} and D_{t_n} in Lemma C.4. We first bound D_{t_n} using the following lemma:

Lemma C.5. *For any $n \in [0 : N - 1]$, we have that*

$$\mathbb{E}[D_{t_n}] \lesssim d\epsilon^2 h_n.$$

Proof. For any $n \in [0 : N - 2]$, we have $T - t_{n+1} \gtrsim \mathcal{O}(1)$ and thus by [111, Corollary 1] that

$$\mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_{N-1} + \tau_{n,m}}(\tilde{\mathbf{x}}_{t_{N-1} + \tau_{n,m}}) - \nabla \log \tilde{p}_{t_{N-1} + \tau'}(\tilde{\mathbf{x}}_{t_{N-1} + \tau'}) \right\|^2 \right] \lesssim d\epsilon_{n,m}^2,$$

for any $\tau' \in [\tau_{n,m}, \tau_{n,m+1}]$, and thus

$$\begin{aligned} \mathbb{E}[D_{t_n}] &= \int_0^{h_n} \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_n + g_n(\tau')}(\tilde{\mathbf{x}}_{t_n + g_n(\tau')}) - \nabla \log \tilde{p}_{t_n + \tau'}(\tilde{\mathbf{x}}_{t_n + \tau'}) \right\|^2 \right] d\tau' \\ &= \sum_{m=0}^{M_n} \int_{\tau_{n,m}}^{\tau_{n,m+1}} \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_n + \tau_{n,m}}(\tilde{\mathbf{x}}_{t_n + \tau_{n,m}}) - \nabla \log \tilde{p}_{t_n + \tau'}(\tilde{\mathbf{x}}_{t_n + \tau'}) \right\|^2 \right] d\tau' \\ &\lesssim \sum_{m=0}^{M_n} d\epsilon_{n,m}^2 \epsilon_{n,m} \leq d\epsilon^2 h_n. \end{aligned}$$

For $n = N - 1$, notice that by the step size schedule (cf. Section 3.1.1) and suppose $\epsilon \leq 1/2$, we have

$$\frac{T - \tau}{2} \leq T - g_n(\tau) \leq T - \tau,$$

and then again [111, Corollary 1] states

$$\mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_n + \epsilon_{n,m}}(\tilde{\mathbf{x}}_{t_n + \epsilon_{n,m}}) - \nabla \log \tilde{p}_{t_n + \tau'}(\tilde{\mathbf{x}}_{t_n + \tau'}) \right\|^2 \right] \lesssim \frac{d\epsilon_{n,m}^2}{T - \tau_{n,m}},$$

and thus

$$\begin{aligned} \mathbb{E}[D_{t_{N-1}}] &= \int_0^{h_{N-1}} \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_{N-1} + g_n(\tau')}(\tilde{\mathbf{x}}_{t_{N-1} + g_n(\tau')}) - \nabla \log \tilde{p}_{t_{N-1} + \tau'}(\tilde{\mathbf{x}}_{t_{N-1} + \tau'}) \right\|^2 \right] d\tau \\ &= \sum_{m=0}^{M_{N-1}} \int_{\tau_{n,m}}^{\tau_{n,m+1}} \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{t_{N-1} + \tau_{n,m}}(\tilde{\mathbf{x}}_{t_{N-1} + \tau_{n,m}}) - \nabla \log \tilde{p}_{t_{N-1} + \tau'}(\tilde{\mathbf{x}}_{t_{N-1} + \tau'}) \right\|^2 \right] d\tau' \\ &\lesssim \sum_{m=0}^{M_{N-1}} \frac{d\epsilon_{n,m}^2}{T - \tau_{n,m}} \epsilon_{n,m} \leq \sum_{m=0}^{M_{N-1}} d\epsilon_{n,m}^2 \epsilon \lesssim \int_{\delta_\infty}^{T - t_{N-1}} d\tau d\tau \lesssim d\epsilon^2 h_{N-1}. \end{aligned}$$

\square

Remark C.6. The above lemma is able to achieve a better dependency on ϵ compared to Lemma B.7, because the backward process $(\tilde{\mathbf{x}}_t)_{t \in [0, T]}$ is now a deterministic process in the probability flow ODE formulation, instead of a stochastic process as in the SDE formulation as in Lemma B.7. Thus, intuitively applying Cauchy-Schwarz rather than Itô symmetry gives us a $\mathcal{O}(\epsilon^2)$ -dependency rather than $\mathcal{O}(\epsilon)$ -dependency.

Theorem C.7. Under Assumptions 3.1', 3.2, 3.3, and 3.4, then the distribution \tilde{q}_{t_n, h_n} that the parallelized predictor step generates samples from satisfies the following error bound:

$$W_2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}})^2 \lesssim de^{-K} + h_n^2 \delta_\infty^2 + d\epsilon^2 h_n^2,$$

for $n \in [0 : N - 1]$.

Proof. By the definition of 2-Wasserstein distance, we have for any coupling of $\tilde{\mathbf{y}}_{t_n, h_n}^{(K)}$ and $\tilde{\mathbf{x}}_{t_n + h_n}$,

$$W_2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}})^2 \leq \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, h_n}^{(K)} - \tilde{\mathbf{x}}_{t_n + h_n} \right\|^2 \right],$$

and therefore

$$\begin{aligned} W_2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}})^2 &\leq \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, h_n}^{(K)} - \tilde{\mathbf{x}}_{t_n + h_n} \right\|^2 \right] \leq \sup_{\tau \in [0, h_n]} \mathbb{E} \left[\left\| \tilde{\mathbf{y}}_{t_n, \tau}^{(K)} - \tilde{\mathbf{x}}_{t_n + \tau} \right\|^2 \right] \\ &\leq 3d \left(\frac{h_n^2 e^{h_n + \frac{3}{2}} L_s^2}{2} \right)^K + \frac{e^{h_n + \frac{3}{2}} h_n / 2}{1 - h_n^2 e^{h_n + \frac{3}{2}} L_s^2 / 2} (h_n \delta_\infty^2 + \mathbb{E}[D_{t_n}]) \\ &\lesssim de^{-K} + h_n^2 \delta_\infty^2 + d\epsilon^2 h_n^2, \end{aligned}$$

where for the second to last inequality we used Lemma C.4, the last inequality is due to Lemma C.5 and the assumption $h_n^2 e^{h_n} L_s^2 \ll 1$. \square

C.3 Parallelized Corrector Step

After each predictor step, we run the corrector step for $\mathcal{O}(1)$ time to reduce the error. Particularly, we apply the Parallelized underdamped Langevin dynamics algorithm [130] to the corrector step, which yields $\mathcal{O}(1)$ approximate time complexity compared to the ordinary implementation of the ULMC dynamics as in [111]. In the following, we will drop the dependency on ω for notational simplicity, and we refer readers to Appendix A.2 and B.2 to review the change of measure arguments and the application of Girsanov's theorem A.4. We will also use a general notation $*^\dagger$ to distinguish the time in the backward process and the inner time in the corrector step of the n -th block.

We first define the true underdamped Langevin dynamics $(\mathbf{u}_{t_n, t^\dagger}, \mathbf{v}_{t_n, t^\dagger})_{t \geq 0}$:

$$\begin{cases} d\mathbf{u}_{t_n, t^\dagger} = \mathbf{v}_{t_n, t^\dagger} dt^\dagger \\ d\mathbf{v}_{t_n, t^\dagger} = -\gamma \mathbf{v}_{t_n, t^\dagger} dt^\dagger - \nabla \log \tilde{p}_{t_{n+1}}(\mathbf{u}_{t_n, t^\dagger}) dt^\dagger + \sqrt{2\gamma} d\mathbf{b}_{t_n, t^\dagger}, \end{cases} \quad (\text{C.10})$$

with initial condition $\mathbf{u}_{t_n, 0} \equiv \tilde{\mathbf{y}}_{t_n, h_n}^{(K^\dagger)}$ from the predictor step and $\mathbf{v}_{t_n, 0} \sim \mathcal{N}(0, \mathbf{I}_d)$, where $(\mathbf{b}_{t_n, t^\dagger})_{t \geq 0}$ is a Wiener process. We may also write the system of SDEs above in the following matrix form:

$$d \begin{bmatrix} \mathbf{u}_{t_n, t^\dagger} \\ \mathbf{v}_{t_n, t^\dagger} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{I}_d \\ \mathbf{0} & -\gamma \mathbf{I}_d \end{bmatrix} \begin{bmatrix} \mathbf{u}_{t_n, t^\dagger} \\ \mathbf{v}_{t_n, t^\dagger} \end{bmatrix} - \begin{bmatrix} \mathbf{0} \\ \nabla \log \tilde{p}_{t_{n+1}}(\mathbf{u}_{t_n, t^\dagger}) \end{bmatrix} dt^\dagger + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} d \begin{bmatrix} \mathbf{b}'_{t_n, t^\dagger} \\ \mathbf{b}_{t_n, t^\dagger} \end{bmatrix}.$$

We run this underdamped Langevin dynamics until the pre-determined time horizon T^\dagger . We also define the joint probability distribution of $(\mathbf{u}_{t_n, t^\dagger}, \mathbf{v}_{t_n, t^\dagger})$ at time t as $\pi_{t_n, t^\dagger}(\mathbf{u}_{t_n, t^\dagger}, \mathbf{v}_{t_n, t^\dagger})$ and its marginal on $\mathbf{u}_{t_n, t^\dagger}$ as $\pi_{t_n, t^\dagger}^{\mathbf{u}}(\mathbf{u}_{t_n, t^\dagger})$.

Similar to the parallelizing strategy in Section 3.1.1, we discretize the time interval $[0, T^\dagger]$ into N^\dagger blocks with length $h^\dagger = T^\dagger / N^\dagger$. Within the n -th block, we further divide the block $[n^\dagger h^\dagger, (n+1)h^\dagger]$ into M^\dagger steps, each with step size $\epsilon^\dagger = h^\dagger / M^\dagger$.

Definition C.8 (Auxiliary corrector process). For any $n^\dagger \in [0 : N^\dagger - 1]$, we define the auxiliary corrector process $(\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)})_{\tau^\dagger \in [0, h^\dagger]}$ as the solution to the following SDE recursively for $k^\dagger \in [0 : K^\dagger - 1]$:

$$\begin{cases} d\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} = \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} d\tau^\dagger, \\ d\hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} = -\gamma \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} d\tau^\dagger - \mathbf{s}_{t_{n+1}}(\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, g_n(\tau^\dagger)}^{(k^\dagger)}) d\tau^\dagger + \sqrt{2\gamma} d\mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^\dagger} \end{cases} \quad (\text{C.11})$$

with the initial condition

$$\begin{cases} \hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)} \equiv \hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger} & \text{for } \tau^\dagger \in [0, h^\dagger], \text{ and } \begin{cases} \hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \equiv \hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger} \\ \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \equiv \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, 0} \end{cases} \text{ for } k \in [1 : K^\dagger], \\ \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)} \equiv \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger} \end{cases} \quad (\text{C.12})$$

where

$$\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger} := \hat{\mathbf{u}}_{t_n, (n^\dagger - 1)h^\dagger, h^\dagger}^{(K^\dagger)}, \quad \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger} := \hat{\mathbf{v}}_{t_n, (n^\dagger - 1)h^\dagger, h^\dagger}^{(K^\dagger)}$$

for $n^\dagger \in [1 : N^\dagger - 1]$, and

$$\hat{\mathbf{u}}_{t_n, 0} = \mathbf{y}_{t_n, h_n}^{(K)}, \quad \hat{\mathbf{v}}_{t_n, 0} \sim \mathcal{N}(0, \mathbf{I}_d).$$

We define the joint probability distribution of $(\hat{\mathbf{u}}_{t_n, t^\dagger}, \hat{\mathbf{v}}_{t_n, t^\dagger})$ at time t as $\hat{\pi}_{t_n, t^\dagger}(\hat{\mathbf{u}}_{t_n, t^\dagger}, \hat{\mathbf{v}}_{t_n, t^\dagger})$ and its marginal on $\hat{\mathbf{u}}_{t_n, t^\dagger}$ as $\hat{\pi}_{t_n, t^\dagger}^{\hat{\mathbf{u}}}(\hat{\mathbf{u}}_{t_n, t^\dagger})$. We will also denote the resulting probability distribution of $\hat{\pi}_{t_n, T^\dagger}^{\hat{\mathbf{u}}}$ as $\hat{q}_{t_{n+1}}$.

Lemma C.9 (Equivalence between (C.2) and (C.11)). For any $n^\dagger \in [0 : N^\dagger - 1]$, the update rule in Algorithm 2 is equivalent to the exact solution of the auxiliary process (C.11) for any $k^\dagger \in [0 : K^\dagger - 1]$ and $\tau^\dagger \in [0, h^\dagger]$.

Proof. Without loss of generality, we will prove the lemma for $m^\dagger = M^\dagger$. The proof for $m^\dagger \in [0 : M^\dagger - 1]$ can be done similarly.

We first rewrite (C.2) into the matrix form:

$$\begin{aligned} d \begin{bmatrix} \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \\ \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \end{bmatrix} &= \begin{bmatrix} \mathbf{0} & \mathbf{I}_d \\ \mathbf{0} & -\gamma \mathbf{I}_d \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \\ \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \end{bmatrix} - \begin{bmatrix} \mathbf{0} \\ \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, g_n(\tau^\dagger)}^{(k^\dagger)}) \end{bmatrix} d\tau^\dagger \\ &+ \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} d \begin{bmatrix} \mathbf{b}'_{t_n, n^\dagger h^\dagger + \tau^\dagger} \\ \mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^\dagger} \end{bmatrix}. \end{aligned} \quad (\text{C.13})$$

Define the time-dependent matrix $\mathbf{G}(\cdot)$ as

$$\mathbf{G}(t^\dagger) := \begin{bmatrix} \mathbf{I}_d & \frac{1-e^{-\gamma t^\dagger}}{\gamma} \mathbf{I}_d \\ \mathbf{0} & e^{-\gamma t^\dagger} \mathbf{I}_d \end{bmatrix} = \exp \left(\begin{pmatrix} \mathbf{0} & \mathbf{I}_d \\ \mathbf{0} & -\gamma \mathbf{I}_d \end{pmatrix} t^\dagger \right), \quad (\text{C.14})$$

satisfying that

$$\frac{d}{dt^\dagger} \mathbf{G}(t^\dagger) = \begin{bmatrix} \mathbf{0} & \mathbf{I}_d \\ \mathbf{0} & -\gamma \mathbf{I}_d \end{bmatrix} \mathbf{G}(t^\dagger) = \mathbf{G}(t^\dagger) \begin{bmatrix} \mathbf{0} & \mathbf{I}_d \\ \mathbf{0} & -\gamma \mathbf{I}_d \end{bmatrix}.$$

Now we multiply $\mathbf{G}(-\tau^\dagger)$ on both sides of (C.13) to obtain:

$$\begin{aligned} d \left(\mathbf{G}(-\tau^\dagger) \begin{bmatrix} \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \\ \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \end{bmatrix} \right) &= -\mathbf{G}(-\tau^\dagger) \begin{bmatrix} \mathbf{0} \\ \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, g_n(\tau^\dagger)}^{(k^\dagger)}) \end{bmatrix} d\tau^\dagger \\ &+ \mathbf{G}(-\tau^\dagger) \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} d \begin{bmatrix} \mathbf{b}'_{t_n, n^\dagger h^\dagger + \tau^\dagger} \\ \mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^\dagger} \end{bmatrix}. \end{aligned}$$

Integrating on both sides from 0 to h^\dagger and multiplying $G(h^\dagger)$ on both sides, we have

$$\begin{aligned}
& \begin{bmatrix} \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau}^{(k^\dagger)} \\ \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau}^{(k^\dagger)} \end{bmatrix} - G(h^\dagger) \begin{bmatrix} \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, 0}^{(k^\dagger)} \\ \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, 0}^{(k^\dagger)} \end{bmatrix} \\
&= - \int_0^{h^\dagger} G(h^\dagger - \tau^{\dagger'}) \begin{bmatrix} \mathbf{0} \\ \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, g(\tau^{\dagger'})}^{(k^\dagger)}) \end{bmatrix} d\tau^{\dagger'} \\
&+ \int_0^{h^\dagger} G(h^\dagger - \tau^{\dagger'}) \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} d \begin{bmatrix} \mathbf{b}'_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} \\ \mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} \end{bmatrix} \\
&= - \sum_{m^\dagger=0}^{M^\dagger-1} \int_{m^\dagger \epsilon^\dagger}^{(m^\dagger+1)\epsilon^\dagger} G(h^\dagger - \tau^{\dagger'}) d\tau^{\dagger'} \begin{bmatrix} \mathbf{0} \\ \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, m^\dagger \epsilon^\dagger}^{(k^\dagger)}) \end{bmatrix} \\
&+ \sum_{m^\dagger=0}^{M^\dagger-1} \int_{m^\dagger \epsilon^\dagger}^{(m^\dagger+1)\epsilon^\dagger} G(h^\dagger - \tau^{\dagger'}) \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} d \begin{bmatrix} \mathbf{b}'_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} \\ \mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} \end{bmatrix} \\
&= - \sum_{m^\dagger=0}^{M^\dagger-1} (G(\epsilon^\dagger) - \mathbf{I}_d) G((M^\dagger - m^\dagger - 1)\epsilon^\dagger) \begin{bmatrix} \mathbf{0} \\ \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, m^\dagger \epsilon^\dagger}^{(k^\dagger)}) \end{bmatrix} \\
&+ \sum_{m^\dagger=0}^{M^\dagger-1} \int_{m^\dagger \epsilon^\dagger}^{(m^\dagger+1)\epsilon^\dagger} G(h^\dagger - \tau^{\dagger'}) \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} d \begin{bmatrix} \mathbf{b}'_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} \\ \mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} \end{bmatrix}.
\end{aligned}$$

By Itô isometry, we have

$$\begin{aligned}
& \int_{m^\dagger \epsilon^\dagger}^{(m^\dagger+1)\epsilon^\dagger} G(h^\dagger - \tau^{\dagger'}) \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} d \begin{bmatrix} \mathbf{b}'_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} \\ \mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} \end{bmatrix} \\
&\sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} G((M^\dagger - m^\dagger - 1)\epsilon^\dagger)^\top (G(\epsilon^\dagger) - \mathbf{I}_d)^\top \right. \\
&\quad \left. (G(\epsilon^\dagger) - \mathbf{I}_d) G((M^\dagger - m^\dagger - 1)\epsilon^\dagger) \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sqrt{2\gamma} \mathbf{I}_d \end{bmatrix} \right) \\
&\sim \left[\mathcal{N} \left(\mathbf{0}, 2\gamma(1 + \gamma^{-2})(1 - e^{-\gamma\epsilon^\dagger})^2 e^{-2\gamma(M^\dagger - m^\dagger + 1)\epsilon^\dagger} \mathbf{I}_d \right) \right],
\end{aligned}$$

as desired \square

Definition C.10 (Interpolating corrector process). *For any $n^\dagger \in [0 : N^\dagger - 1]$, we define the interpolating corrector process $(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)})_{\tau^\dagger \in [0, h^\dagger]}$ as the solution to the following SDE recursively for $k^\dagger \in [0 : K^\dagger - 1]$:*

$$\begin{cases} d\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} = \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} d\tau^\dagger, \\ d\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} = -\gamma \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} d\tau^\dagger - \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, g_n(\tau^\dagger)}^{(k^\dagger)}) d\tau^\dagger + \sqrt{2\gamma} d\mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^\dagger} \end{cases} \quad (\text{C.15})$$

with the initial condition

$$\begin{cases} \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)} \equiv \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger} & \text{for } \tau^\dagger \in [0, h^\dagger], \text{ and } \\ \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)} \equiv \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger} & \text{for } \tau^\dagger \in [0, h^\dagger], \end{cases} \quad \text{for } k \in [1 : K^\dagger], \quad (\text{C.16})$$

where

$$\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger} := \tilde{\mathbf{u}}_{(n^\dagger-1)h^\dagger, h^\dagger}^{(K^\dagger)}, \quad \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger} := \tilde{\mathbf{v}}_{(n^\dagger-1)h^\dagger, h^\dagger}^{(K^\dagger)}$$

for $n^\dagger \in [1 : N^\dagger - 1]$, and

$$\tilde{\mathbf{u}}_{t_n, 0} = \tilde{\mathbf{y}}_{t_n, h_n}^{(K)}, \quad \tilde{\mathbf{v}}_{t_n, 0} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d).$$

We define the joint probability distribution of $(\tilde{\mathbf{u}}_{t_n, t^\dagger}, \tilde{\mathbf{v}}_{t_n, t^\dagger})$ at time t as $\tilde{\pi}_{t_n, t^\dagger}(\tilde{\mathbf{u}}_{t_n, t^\dagger}, \tilde{\mathbf{v}}_{t_n, t^\dagger})$ and its marginal on $\tilde{\mathbf{u}}_{t_n, t^\dagger}$ as $\tilde{\pi}_{t_n, t^\dagger}^{\tilde{\mathbf{u}}}(\tilde{\mathbf{u}}_{t_n, t^\dagger})$.

We invoke Girsanov's theorem (Theorem A.4) again by the following procedure

1. Setting (A.2) as the auxiliary process (C.15) at iteration K^\dagger , where $\mathbf{b}_{t_n, t^\dagger}(\omega)$ is a Wiener process under the measure Q ;
2. Defining another process $\tilde{\mathbf{b}}_{t_n, n^\dagger h^\dagger + \tau^\dagger}$ governed by the following SDE:

$$d\tilde{\mathbf{b}}_{t_n, n^\dagger h^\dagger + \tau^\dagger} = d\mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^\dagger} - \phi_{t_n, n^\dagger h^\dagger}(\tau^\dagger) d\tau^\dagger, \quad (\text{C.17})$$

where

$$\phi_{t_n, n^\dagger h^\dagger}(\tau^\dagger) = \frac{1}{\sqrt{2\gamma}} \left(\mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}) - \nabla \log \tilde{p}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}) \right) \quad (\text{C.18})$$

and computing the Radon-Nikodym derivative of the measure P with respect to Q as

$$\frac{dP}{dQ} = \exp \left(\int_0^{h^\dagger} \phi_{t_n, n^\dagger h^\dagger}(\tau^\dagger)^\top d\mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^\dagger} - \frac{1}{2} \int_0^{h^\dagger} \|\phi_{t_n, n^\dagger h^\dagger}(\tau^\dagger)\|^2 d\tau^\dagger \right); \quad (\text{C.19})$$

3. Concluding that (C.15) at iteration K^\dagger under the measure Q satisfies the following SDE:

$$\begin{cases} d\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} = \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} d\tau^\dagger \\ d\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} = -\gamma \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} d\tau^\dagger - \nabla \log \tilde{p}_{t_{n+1}}(\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}) d\tau^\dagger + \sqrt{2\gamma} d\tilde{\mathbf{b}}_{t_n, n^\dagger h^\dagger + \tau^\dagger}, \end{cases} \quad (\text{C.20})$$

with $(\tilde{\mathbf{b}}_{t_n, n^\dagger h^\dagger + \tau^\dagger})_{\tau^\dagger \geq 0}$ being a Wiener process under the measure P . If we replace $(\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}, \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)})$ by $(\mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}, \mathbf{v}_{t_n, n^\dagger h^\dagger + \tau^\dagger})$, one should notice (C.20) is immediately the original backward SDE (C.10) with the true score function on $t \in [n^\dagger h^\dagger, (n+1)h^\dagger]$:

$$\begin{cases} d\mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger} = \mathbf{v}_{t_n, n^\dagger h^\dagger + \tau^\dagger} d\tau^\dagger \\ d\mathbf{v}_{t_n, n^\dagger h^\dagger + \tau^\dagger} = -\gamma \mathbf{v}_{t_n, n^\dagger h^\dagger + \tau^\dagger} d\tau^\dagger - \nabla \log \tilde{p}_{t_{n+1}}(\mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}) d\tau^\dagger + \sqrt{2\gamma} d\tilde{\mathbf{b}}_{t_n, n^\dagger h^\dagger + \tau^\dagger}. \end{cases} \quad (\text{C.21})$$

We further define the joint probability distribution of $(\mathbf{u}_{t_n, t^\dagger}, \mathbf{v}_{t_n, t^\dagger})$ at time t as $\pi_{t_n, t^\dagger}(\mathbf{u}_{t_n, t^\dagger}, \mathbf{v}_{t_n, t^\dagger})$ and its marginal on $\mathbf{u}_{t_n, t^\dagger}$ as $\pi_{t_n, t^\dagger}^{\mathbf{u}}(\mathbf{u}_{t_n, t^\dagger})$.

Remark C.11. The application of Girsanov's theorem A.4 is by writing the system of SDEs in the matrix form.

Definition C.12 (Stationary process). Under the P -measure that is defined by the Radon-Nikodym derivative (C.19), we may define a stationary underdamped Langevin process for $n^\dagger \in [0 : N^\dagger - 1]$ and $\tau^\dagger \in [0, h^\dagger]$ as

$$\begin{cases} d\mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^* = \mathbf{v}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^* d\tau^\dagger, \\ d\mathbf{v}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^* = -\gamma \mathbf{v}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^* d\tau^\dagger - \nabla \log \tilde{p}_{t_{n+1}}(\mathbf{u}_{n^\dagger h^\dagger + \tau^\dagger}^*) d\tau^\dagger + \sqrt{2\gamma} d\tilde{\mathbf{b}}_{t_n, n^\dagger h^\dagger + \tau^\dagger}, \end{cases} \quad (\text{C.22})$$

with the initial condition $\mathbf{u}_{t_n, n^\dagger h^\dagger}^* \sim \tilde{p}_{t_{n+1}}$ and $\mathbf{v}_{t_n, n^\dagger h^\dagger}^* \sim \mathcal{N}(0, \mathbf{I}_d)$. We define the joint probability distribution of $(\mathbf{u}_{t_n, t^\dagger}^*, \mathbf{v}_{t_n, t^\dagger}^*)$ at time t as $\pi_{t_n, t^\dagger}^{\mathbf{u}^*}(\mathbf{u}_{t_n, t^\dagger}^*, \mathbf{v}_{t_n, t^\dagger}^*)$ and its marginal on $\mathbf{u}_{t_n, t^\dagger}^*$ as $\pi_{t_n, t^\dagger}^{\mathbf{u}^*}(\mathbf{u}_{t_n, t^\dagger}^*)$.

Thus, from Corollary A.5, we have that

$$\begin{aligned} & D_{\text{KL}}(\pi_{t_n, n^\dagger h^\dagger} \| \tilde{\pi}_{t_n, n^\dagger h^\dagger}) \\ & \leq D_{\text{KL}}(\pi_{t_n, (n-1)h^\dagger} \| \tilde{\pi}_{t_n, (n-1)h^\dagger}) + \sum_{n=0}^{N^\dagger-1} D_{\text{KL}}(\pi_{t_n, n^\dagger h^\dagger : (n+1)h^\dagger} \| \tilde{\pi}_{t_n, n^\dagger h^\dagger : (n+1)h^\dagger}) \\ & \leq D_{\text{KL}}(\pi_{t_n, (n-1)h^\dagger} \| \tilde{\pi}_{t_n, (n-1)h^\dagger}) \\ & \quad + \frac{1}{4\gamma} \mathbb{E}_P \left[\int_0^{h^\dagger} \left\| \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}) - \nabla \log \tilde{p}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}) \right\|^2 d\tau^\dagger \right]. \end{aligned} \quad (\text{C.23})$$

By triangle inequality, we have

$$\begin{aligned}
& \int_0^{h^\dagger} \left\| \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)}) - \nabla \log \tilde{p}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}) \right\|^2 d\tau^\dagger \\
& \leq 5 \int_0^{h^\dagger} \left\| \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)}) - \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)}) \right\|^2 d\tau^\dagger \\
& + 5 \int_0^{h^\dagger} \left\| \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)}) - \mathbf{s}_{t_{n+1}}(\mathbf{u}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^*) \right\|^2 d\tau^\dagger \\
& + 5 \int_0^{h^\dagger} \left\| \mathbf{s}_{t_{n+1}}(\mathbf{u}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^*) - \nabla \log \tilde{p}_{t_{n+1}}(\mathbf{u}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^*) \right\|^2 d\tau^\dagger \\
& + 5 \int_0^{h^\dagger} \left\| \nabla \log \tilde{p}_{t_{n+1}}(\mathbf{u}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^*) - \nabla \log \tilde{p}_{t_{n+1}}(\mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^*) \right\|^2 d\tau^\dagger \\
& + 5 \int_0^{h^\dagger} \left\| \nabla \log \tilde{p}_{t_{n+1}}(\mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^*) - \nabla \log \tilde{p}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}) \right\|^2 d\tau^\dagger \\
& \leq 5L_s^2 \int_0^{h^\dagger} \left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)} \right\|^2 d\tau^\dagger \\
& + 5 \underbrace{\left(L_s^2 \int_0^{h^\dagger} \left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)} - \mathbf{u}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^* \right\|^2 d\tau^\dagger + L_p^2 \int_0^{h^\dagger} \left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} - \mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^* \right\|^2 d\tau^\dagger \right)}_{:=E_{t_n, n^\dagger h^\dagger}} \\
& + 5h^\dagger \delta_\infty^2 + 5L_p^2 \underbrace{\int_0^{h^\dagger} \left\| \mathbf{u}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^* - \mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^* \right\|^2 d\tau^\dagger}_{:=F_{t_n, n^\dagger h^\dagger}}, \tag{C.24}
\end{aligned}$$

where we used the Lipschitz continuity of the learned score function (Assumption 3.3) and the true score function (Assumption 3.4), and the δ_∞ -accuracy of the learned score function at each time step (Assumption 3.1').

Now we proceed to bound the terms in the error decomposition (C.24). We first bound the $F_{t_n, n^\dagger h^\dagger}$ term by the following lemma:

Lemma C.13. *For any $n \in [0 : N - 1]$ and $\tau^\dagger \in [0, h^\dagger]$, we have*

$$\mathbb{E}_P \left[\left\| \mathbf{u}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^* - \mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^* \right\|^2 \right] \leq d\epsilon^{\dagger 2},$$

and therefore

$$\mathbb{E}_P [F_{t_n, n^\dagger h^\dagger}] \leq dh^\dagger \epsilon^{\dagger 2}.$$

Proof. By the definition of $(\mathbf{u}_{t_n, n^\dagger h^\dagger, \tau}^*, \mathbf{v}_{t_n, n^\dagger h^\dagger, \tau}^*)$ as the stationary underdamped Langevin dynamics (C.22), we have

$$\begin{aligned}
\mathbb{E}_P \left[\left\| \mathbf{u}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^* - \mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^* \right\|^2 \right] &= \mathbb{E}_P \left[\left\| \int_{\lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{\tau^\dagger} \mathbf{v}_{t_n, n^\dagger h^\dagger, \tau'^*}^* d\tau'^* \right\|^2 \right] \\
&\leq \epsilon^\dagger \int_{\lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{\tau^\dagger} \mathbb{E}_P \left[\left\| \mathbf{v}_{t_n, n^\dagger h^\dagger, \tau'^*}^* \right\|^2 \right] d\tau'^* \leq d\epsilon^{\dagger 2},
\end{aligned}$$

where the first inequality follows from Cauchy-Schwarz inequality and the last inequality is by the fact that

$$\mathbf{v}_{t_n, n^\dagger h^\dagger, \tau'^*}^* \sim \mathcal{N}(0, \mathbf{I}_d), \quad \text{for any } \tau'^* \in [0, h^\dagger].$$

Consequently, we have

$$\mathbb{E}_P [F_{t_n, n^\dagger h^\dagger}] = \int_0^{h^\dagger} \mathbb{E}_P \left[\left\| \mathbf{u}_{t_n, n^\dagger h^\dagger + \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^* - \mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^* \right\|^2 \right] d\tau^\dagger \leq dh^\dagger \epsilon^{\dagger 2}.$$

□

The term $E_{t_n, n^\dagger h^\dagger}$ can be bounded with the following lemma:

Lemma C.14. *For any $n^\dagger \in [0 : N^\dagger - 1]$, suppose that $\gamma \lesssim L_p^{-1/2}$ and $T^\dagger \lesssim L_p^{-1/2}$, then we have the following inequality for any $\tau^\dagger \in [0, h^\dagger]$*

$$\mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} - \mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^* \right\|^2 \right] \lesssim W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}),$$

and therefore

$$\mathbb{E}_P [E_{t_n, n^\dagger h^\dagger}] \lesssim h^\dagger (L_s^2 + L_p^2) W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}).$$

Proof. Recall that under the measure P , $\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}$ follows the dynamics of $\mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}$ (C.21) for $\tau^\dagger \in [0, h^\dagger]$, which coincides with that of $\mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^*$. As the only difference between the two processes $\mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}$ and $\mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^*$ is the initial condition, we can invoke Lemma 10 proved in [111] to deduce that

$$\mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} - \mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^* \right\|^2 \right] \lesssim W_2^2(\pi_{t_n, n^\dagger h^\dagger}, \tilde{p}_{t_{n+1}}),$$

where the assumption that $\gamma \lesssim L_p^{-1/2}$ and $T^\dagger \lesssim L_p^{-1/2}$ is required.

Now notice that $\mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^*$ and $\mathbf{u}_{t_n, n^\dagger h^\dagger, \tau^\dagger}$ also follow the same dynamics with the true score function for $\tau^\dagger \in [0, n^\dagger h^\dagger]$, for any coupling of $\mathbf{u}_{t_n, n^\dagger h^\dagger}^*$ and $\mathbf{u}_{t_n, n^\dagger h^\dagger}$, we have

$$\begin{aligned} W_2^2(\pi_{t_n, n^\dagger h^\dagger}, \tilde{p}_{t_{n+1}}) &\leq \mathbb{E} \left[\left\| \mathbf{u}_{t_n, n^\dagger h^\dagger} - \mathbf{u}_{t_n, n^\dagger h^\dagger}^* \right\|^2 \right] \\ &\leq W_2^2(\pi_{t_n, 0}, \tilde{p}_{t_{n+1}}) = W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}), \end{aligned}$$

where the last equality is again by [111, Lemma 10].

Therefore, we have

$$\begin{aligned} &\mathbb{E}_P [E_{t_n, n^\dagger h^\dagger}] \\ &= \int_0^{h^\dagger} \mathbb{E}_P \left[L_s^2 \left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)} - \mathbf{u}_{t_n, n^\dagger h^\dagger + \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^* \right\|^2 + L_p^2 \left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} - \mathbf{u}_{t_n, n^\dagger h^\dagger + \tau^\dagger}^* \right\|^2 \right] d\tau^\dagger \\ &\leq h^\dagger (L_s^2 + L_p^2) W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}). \end{aligned}$$

□

Now, we provide lemmas that are used to bound the first term in (C.24).

Lemma C.15. *For any $n^\dagger \in [0 : N^\dagger - 1]$, we have the following estimate:*

$$\begin{aligned} &\sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)} \right\|^2 \right] \\ &\leq \frac{5L_s^2 h^\dagger e^{(3+\gamma)h^\dagger}}{2\gamma} \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger-1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} \right\|^2 \right] \\ &\quad + \frac{5h^\dagger e^{(3+\gamma)h^\dagger}}{2\gamma} \mathbb{E}_P [E_{t_n, n^\dagger h^\dagger} + h^\dagger \delta_\infty^2 + L_p^2 F_{t_n, n^\dagger h^\dagger}] + h^{\dagger 2} e^{(3+\gamma)h^\dagger} (3\gamma d + M_s^2) + h^\dagger e^{2h^\dagger} d. \end{aligned}$$

Proof. Let $\mu_{t_n, n^\dagger h^\dagger, \tau^\dagger} := \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)}$ and $\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger} := \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(1)} - \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)}$. Then for $k = 0$, we may rewrite (C.15) as follows

$$\begin{cases} d\mu_{t_n, n^\dagger h^\dagger, \tau^\dagger} = \left(\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger} + \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)} \right) d\tau^\dagger \\ d\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger} = -\gamma(\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger} + \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)}) d\tau^\dagger - \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)}) d\tau^\dagger + \sqrt{2\gamma} d\mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^\dagger} \end{cases} \quad (\text{C.25})$$

On the one hand, by using the first equation in (C.25), we may compute the derivative

$$\frac{d}{d\tau^{\dagger'}} \|\mu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}\|^2 = 2\mu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^\top \left(\nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}} + \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^{(0)} \right)$$

and integrate it for $\tau^{\dagger'} \in [0, \tau^\dagger]$, which yields

$$\begin{aligned} \|\mu_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 &= 2 \int_0^{\tau^\dagger} \mu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^\top (\nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}} + \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^{(0)}) d\tau^{\dagger'} \\ &\leq 2 \int_0^{\tau^\dagger} \|\mu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}\|^2 d\tau^{\dagger'} + \int_0^{\tau^\dagger} \|\nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}\|^2 d\tau^{\dagger'} + \int_0^{\tau^\dagger} \|\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^{(0)}\|^2 d\tau^{\dagger'}. \end{aligned}$$

Applying Gronwall's inequality, we have

$$\|\mu_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 \leq e^{2\tau^\dagger} \left(\int_0^{\tau^\dagger} \|\nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}\|^2 d\tau^{\dagger'} + \int_0^{\tau^\dagger} \|\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^{(0)}\|^2 d\tau^{\dagger'} \right).$$

We then take expectation with respect to the path measure P and then the supremum with respect to $\tau^\dagger \in [0, h^\dagger]$, implying that

$$\begin{aligned} &\sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\|\mu_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 \right] \\ &\leq \sup_{\tau^\dagger \in [0, h^\dagger]} \left(e^{2\tau^\dagger} \int_0^{\tau^\dagger} \mathbb{E}_P \left[\|\nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}\|^2 \right] d\tau^{\dagger'} + e^{2\tau^\dagger} \int_0^{\tau^\dagger} \mathbb{E}_P \left[\|\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^{(0)}\|^2 \right] d\tau^{\dagger'} \right) \\ &\leq h^\dagger e^{2h^\dagger} \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\|\nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}\|^2 \right] + h^\dagger e^{2h^\dagger} d. \end{aligned} \quad (\text{C.26})$$

On the other hand, by applying Itô's lemma and plugging in the expression of $\mathbf{b}_{t_n, n^\dagger h^\dagger + \tau^\dagger}$ given by (C.17), we have

$$\begin{aligned} &d\|\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 \\ &= - \left[2\gamma \|\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 + 2\gamma \nu_{t_n, n^\dagger h^\dagger, \tau^\dagger}^\top \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)} + 2\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger}^\top \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)}) - 2\gamma d \right] d\tau^\dagger \\ &\quad + 2\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger}^\top \sqrt{2\gamma} (d\tilde{\mathbf{b}}_{t_n, n^\dagger h^\dagger + \tau^\dagger} + \phi_{t_n, n^\dagger h^\dagger}(\tau^\dagger) d\tau^\dagger), \end{aligned} \quad (\text{C.27})$$

Then similarly, we may compute the derivative of $\|\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2$, integrate it for $\tau^\dagger \in [0, h^\dagger]$, and take the supremum with respect to τ^\dagger to obtain

$$\begin{aligned} &\mathbb{E}_P \left[\|\nu_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 \right] \\ &= \mathbb{E}_P \left[- \int_0^{\tau^\dagger} \left(2\gamma \|\nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}\|^2 + 2\gamma \nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^\top \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^{(0)} - 2\gamma d \right) d\tau^{\dagger'} \right] \\ &\quad + \mathbb{E}_P \left[- \int_0^{\tau^\dagger} 2\nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^\top \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^{(0)}) d\tau^{\dagger'} \right] \\ &\quad + 2\sqrt{2\gamma} \mathbb{E}_P \left[\int_0^{\tau^\dagger} \nu_{t_n, n^\dagger h^\dagger, \tau^{\dagger'}}^\top (d\tilde{\mathbf{b}}_{t_n, n^\dagger h^\dagger + \tau^{\dagger'}} + \phi_{t_n, n^\dagger h^\dagger}(\tau^{\dagger'}) d\tau^{\dagger'}) \right]. \end{aligned}$$

By Itô's lemma, this equals to

$$\begin{aligned} & \mathbb{E}_P \left[\|\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 \right] \\ &= \mathbb{E}_P \left[- \int_0^{\tau^\dagger} \left(2\gamma \|\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}\|^2 + 2\gamma \boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}^\top \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}^{(0)} - 2\gamma d \right) d\tau^\dagger' \right] \\ &+ \mathbb{E}_P \left[- \int_0^{\tau^\dagger} 2\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}^\top \mathbf{s}_{t_{n+1}} \left(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}^{(0)} \right) + 2\sqrt{2\gamma} \boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}^\top \boldsymbol{\phi}_{t_n, n^\dagger h^\dagger}(\tau^\dagger') d\tau^\dagger' \right]. \end{aligned}$$

Applying AM-GM gives

$$\begin{aligned} & \mathbb{E}_P \left[\|\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 \right] \\ &\leq \int_0^{\tau^\dagger} \mathbb{E}_P \left[(1+\gamma) \|\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}\|^2 + \|\boldsymbol{\phi}_{t_n, n^\dagger h^\dagger}(\tau^\dagger')\|^2 \right] d\tau^\dagger' \\ &+ \int_0^{\tau^\dagger} \mathbb{E}_P \left[\gamma \|\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}^{(0)}\|^2 + \|\mathbf{s}_{t_{n+1}} \left(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}^{(0)} \right)\|^2 + 2\gamma d \right] d\tau^\dagger' \\ &\leq \int_0^{\tau^\dagger} \mathbb{E}_P \left[(1+\gamma) \|\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}\|^2 + \|\boldsymbol{\phi}_{t_n, n^\dagger h^\dagger}(\tau^\dagger')\|^2 \right] d\tau^\dagger' + \left(\gamma \mathbb{E} \left[\|\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, 0}^{(0)}\|^2 \right] + M_s^2 + 2\gamma d \right) \tau^\dagger \\ &= (1+\gamma) \int_0^{\tau^\dagger} \mathbb{E}_P \left[\|\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger'}\|^2 \right] d\tau^\dagger' + \int_0^{\tau^\dagger} \mathbb{E}_P \left[\|\boldsymbol{\phi}_{t_n, n^\dagger h^\dagger}(\tau^\dagger')\|^2 \right] d\tau^\dagger' + \tau^\dagger (3\gamma d + M_s^2), \end{aligned}$$

where in the last equality, we used the initialization of the auxiliary corrector process $\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, 0}^{(0)} \sim \mathcal{N}(0, \mathbf{I}_d)$.

Again, we apply Gronwall's inequality to the above inequality and take the supremum with respect to $\tau^\dagger \in [0, h^\dagger]$ to obtain

$$\begin{aligned} & \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\|\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 \right] \\ &\leq e^{(1+\gamma)h^\dagger} \int_0^{h^\dagger} \mathbb{E}_P \left[\|\boldsymbol{\phi}_{t_n, n^\dagger h^\dagger}(\tau^\dagger)\|^2 \right] d\tau^\dagger + h^\dagger e^{(1+\gamma)h^\dagger} (3\gamma d + M_s^2) \\ &\leq \frac{e^{(1+\gamma)h^\dagger}}{2\gamma} \mathbb{E}_P \left[\int_0^{h^\dagger} \left\| \mathbf{s}_{t_{n+1}} \left(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)} \right) - \nabla \log \tilde{p}_{t_{n+1}} \left(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} \right) \right\|^2 d\tau^\dagger \right] \\ &\quad + h^\dagger e^{(1+\gamma)h^\dagger} (3\gamma d + M_s^2), \end{aligned} \tag{C.28}$$

and for the difference term within the expectation, we decompose it again by the triangle inequality in (C.24), *i.e.*

$$\begin{aligned} & \int_0^{h^\dagger} \left\| \mathbf{s}_{t_{n+1}} \left(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)} \right) - \nabla \log \tilde{p}_{t_{n+1}} \left(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} \right) \right\|^2 d\tau^\dagger \\ &\leq 5L_s^2 \int_0^{h^\dagger} \left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)} \right\|^2 d\tau^\dagger + 5E_{t_n, n^\dagger h^\dagger} + 5h^\dagger \delta_\infty^2 + 5L_p^2 F_{t_n, n^\dagger h^\dagger}, \end{aligned}$$

to obtain that

$$\begin{aligned} & \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\|\boldsymbol{\nu}_{t_n, n^\dagger h^\dagger, \tau^\dagger}\|^2 \right] \\ &\leq \frac{5L_s^2 e^{(1+\gamma)h^\dagger}}{2\gamma} \mathbb{E}_P \left[\int_0^{h^\dagger} \left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)} \right\|^2 d\tau^\dagger \right] \\ &+ \frac{5e^{(1+\gamma)h^\dagger}}{2\gamma} \mathbb{E}_P \left[E_{t_n, n^\dagger h^\dagger} + h^\dagger \delta_\infty^2 + L_p^2 F_{t_n, n^\dagger h^\dagger} \right] + h^\dagger e^{(1+\gamma)h^\dagger} (3\gamma d + M_s^2) \\ &\leq \frac{5L_s^2 e^{(1+\gamma)h^\dagger}}{2\gamma} h^\dagger \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)} \right\|^2 \right] \\ &+ \frac{5e^{(1+\gamma)h^\dagger}}{2\gamma} \mathbb{E}_P \left[E_{t_n, n^\dagger h^\dagger} + h^\dagger \delta_\infty^2 + L_p^2 F_{t_n, n^\dagger h^\dagger} \right] + h^\dagger e^{(1+\gamma)h^\dagger} (3\gamma d + M_s^2), \end{aligned}$$

substituting which into (C.26) completes our proof of this Lemma. \square

Lemma C.16 (Exponential convergence of Picard iteration in the corrector step of PIADM-ODE). *For any $n^\dagger \in [0, N^\dagger - 1]$, then the two ending terms $\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}$ and $\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)}$ of the sequence $\{\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)}\}_{k^\dagger \in [0:K^\dagger-1]}$ satisfy the following exponential convergence rate*

$$\begin{aligned} & \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger-1)} \right\|^2 \right] \\ & \leq C_{K^\dagger} \left(\frac{5h^\dagger e^{(3+\gamma)h^\dagger}}{2\gamma} \mathbb{E}_P [E_{t_n, n^\dagger h^\dagger} + h^\dagger \delta_\infty^2 + L_p^2 F_{t_n, n^\dagger h^\dagger}] + h^{\dagger 2} e^{(3+\gamma)h^\dagger} (3\gamma d + M_s^2) + h^\dagger e^{2h^\dagger} d \right), \end{aligned} \quad (\text{C.29})$$

where the coefficient

$$C_{K^\dagger} = \left(\frac{L_s^2 h^{\dagger 2} e^{h^\dagger}}{2\gamma} \right)^{K^\dagger-1} \bigg/ \left(1 - \frac{5L_s^2 h^\dagger e^{(3+\gamma)h^\dagger}}{2\gamma} \left(\frac{L_s^2 h^{\dagger 2} e^{h^\dagger}}{2\gamma} \right)^{K^\dagger-1} \right).$$

Proof. We subtract the dynamics of $\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)}$ and $\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k)}$ in (C.15) to obtain

$$d \left(\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right) = \left(\tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right) d\tau^\dagger.$$

Then, we use the formula above to compute the derivative

$$\frac{d}{d\tau^{\dagger'}} \left\| \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right\|^2 = 2 \left(\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right)^\top \left(\tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right)$$

and integrate for $\tau^{\dagger'} \in [0, \tau^\dagger]$ to obtain

$$\begin{aligned} & \left\| \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right\|^2 \\ & = 2 \int_0^{\tau^\dagger} \left(\tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right)^\top \left(\tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right) d\tau^{\dagger'} \\ & \leq \int_0^{\tau^\dagger} \left\| \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right\|^2 d\tau^{\dagger'} + \int_0^{\tau^\dagger} \left\| \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right\|^2 d\tau^{\dagger'} \end{aligned}$$

Applying Grönwall's inequality gives us that

$$\left\| \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right\|^2 \leq e^{\tau^\dagger} \int_0^{\tau^\dagger} \left\| \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right\|^2 d\tau^{\dagger'}$$

and taking the supremum with respect to $\tau^\dagger \in [0, h^\dagger]$ on both sides above implies

$$\sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{u}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right\|^2 \right] \leq h^\dagger e^{h^\dagger} \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^{\dagger'}}^{(k^\dagger)} \right\|^2 \right]. \quad (\text{C.30})$$

We then apply a similar argument for $\tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)}$ as well

$$\begin{aligned} & d \left(\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right) \\ & = -\gamma \left(\tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right) d\tau^\dagger - \left(\mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}) - \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}) \right) d\tau^\dagger, \end{aligned}$$

integrate which for $\tau^\dagger \in [0, \tau^\dagger]$ to obtain

$$\begin{aligned}
& \left\| \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right\|^2 \\
&= - \int_0^{\tau^\dagger} 2\gamma \left\| \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right\|^2 d\tau^\dagger \\
&\quad - 2 \int_0^{\tau^\dagger} \left(\tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right)^\top \left(\mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(k^\dagger)}) - \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(k-1)}) \right) d\tau^\dagger \\
&\leq \frac{1}{2\gamma} \int_0^{\tau^\dagger} \left\| \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(k^\dagger)}) - \mathbf{s}_{t_{n+1}}(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(k-1)}) \right\|^2 d\tau^\dagger \\
&\leq \frac{L_s^2}{2\gamma} \int_0^{\tau^\dagger} \left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(k^\dagger)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(k-1)} \right\|^2 d\tau^\dagger.
\end{aligned}$$

And then taking the supremum with respect to $\tau^\dagger \in [0, h^\dagger]$ on both sides above implies

$$\sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k+1)} - \tilde{\mathbf{v}}_{n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} \right\|^2 \right] \leq \frac{h^\dagger L_s^2}{2\gamma} \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k^\dagger)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(k-1)} \right\|^2 \right] \quad (\text{C.31})$$

Substituting (C.31) into (C.30) and iterating over $k \in [1 : K^\dagger - 1]$, we obtain that

$$\begin{aligned}
& \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger-1)} \right\|^2 \right] \leq \frac{L_s^2 h^{\dagger 2} e^{h^\dagger}}{2\gamma} \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger-1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger-2)} \right\|^2 \right] \\
&\leq \left(\frac{L_s^2 h^{\dagger 2} e^{h^\dagger}}{2\gamma} \right)^{K^\dagger-1} \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(0)} \right\|^2 \right] \\
&\leq \left(\frac{L_s^2 h^{\dagger 2} e^{h^\dagger}}{2\gamma} \right)^{K^\dagger-1} \frac{5h^\dagger e^{(3+\gamma)h^\dagger}}{2\gamma} \mathbb{E}_P [E_{t_n, n^\dagger h^\dagger} + h^\dagger \delta_\infty^2 + L_p^2 F_{t_n, n^\dagger h^\dagger}] \\
&\quad + \left(\frac{L_s^2 h^{\dagger 2} e^{h^\dagger}}{2\gamma} \right)^{K^\dagger-1} \left(h^{\dagger 2} e^{(3+\gamma)h^\dagger} (3\gamma d + M_s^2) + h^\dagger e^{2h^\dagger} d \right) \\
&\quad + \left(\frac{L_s^2 h^{\dagger 2} e^{h^\dagger}}{2\gamma} \right)^{K^\dagger-1} \frac{5L_s^2 h^\dagger e^{(3+\gamma)h^\dagger}}{2\gamma} \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger-1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} \right\|^2 \right],
\end{aligned}$$

where we plug in the results from Lemma C.15 in the last inequality. Rearranging the inequality above completes our proof. \square

Theorem C.17. Under Assumptions 3.1', 3.2, 3.3, and 3.4, given the following choices of the order of the parameters

$$\begin{aligned}
T^\dagger &= \mathcal{O}(1), \quad N^\dagger = \mathcal{O}(1), \quad h^\dagger = \Theta(1) \\
M^\dagger &= \Theta(d^{1/2} \delta^{-1}), \quad \epsilon^\dagger = \Theta(d^{-1/2} \delta), \quad K^\dagger = \mathcal{O}(\log(d\delta^{-2}))
\end{aligned}$$

and let

$$\frac{L_s^2 h^{\dagger 2} e^{h^\dagger}}{2\gamma} \ll 1, \quad \gamma \lesssim L_p^{-1/2}, \quad T^\dagger \lesssim L_p^{-1/2} \wedge L_s^{-1/2}, \quad \delta_\infty \lesssim \delta$$

then the distribution $\tilde{\pi}_{t_n, T^\dagger}$ satisfies the following error bound:

$$\begin{aligned}
D_{\text{KL}}(\pi_{t_n, T^\dagger} \| \tilde{\pi}_{t_n, T^\dagger}) &\lesssim T^\dagger W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + T^\dagger \delta_\infty^2 + d T^\dagger \epsilon^{\dagger 2} + e^{-K^\dagger} T^\dagger h^\dagger d \\
&\lesssim W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + \delta^2,
\end{aligned}$$

with a total of $K^\dagger N^\dagger = \mathcal{O}(\log(d\delta^{-2}))$ approximate time complexity and $M = \Theta(d^{1/2} \delta^{-2})$ space complexity for parallelizable δ -accurate score function computations.

Proof. Now, we continue the computation by plugging the decomposition in (C.24) and all the error bounds derived above into the equation. First for the last term in (C.23)

$$\begin{aligned}
& \mathbb{E}_P \left[\int_0^{h^\dagger} \left\| \mathbf{s}_{t_{n+1}} \left(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)} - \nabla \log \tilde{p}_{t_{n+1}} \left(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}^{(K^\dagger)} \right) \right\|^2 d\tau^\dagger \right] \\
& \leq 5L_s^2 h^\dagger \sup_{\tau^\dagger \in [0, h^\dagger]} \mathbb{E}_P \left[\left\| \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger-1)} - \tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \lfloor \frac{\tau^\dagger}{\epsilon^\dagger} \rfloor \epsilon^\dagger}^{(K^\dagger)} \right\|^2 \right] + 5\mathbb{E}_P [E_{t_n, n^\dagger h^\dagger} + h^\dagger \delta_\infty^2 + L_p^2 F_{t_n, n^\dagger h^\dagger}] \\
& \leq 5 \left(1 + L_s^2 h^\dagger C_{K^\dagger} \frac{5h^\dagger e^{(3+\gamma)h^\dagger}}{2\gamma} \right) \mathbb{E}_P [E_{t_n, n^\dagger h^\dagger} + h^\dagger \delta_\infty^2 + L_p^2 F_{t_n, n^\dagger h^\dagger}] \\
& \quad + 5L_s^2 h^\dagger C_{K^\dagger} \left(h^{\dagger 2} e^{(3+\gamma)h^\dagger} (3\gamma d + M_s^2) + h^\dagger e^{2h^\dagger} d \right),
\end{aligned}$$

where the last inequality is by Lemma C.16. We further substitute Lemma C.14 and C.13 into (C.23) to obtain

$$\begin{aligned}
& D_{\text{KL}}(\pi_{t_n, n^\dagger h^\dagger} \| \tilde{\pi}_{t_n, n^\dagger h^\dagger}) \\
& \leq D_{\text{KL}}(\pi_{t_n, (n-1)h^\dagger} \| \tilde{\pi}_{t_n, (n-1)h^\dagger}) + 5 \frac{2\gamma + L_s^2 C_{K^\dagger} 5h^{\dagger 2} e^{(3+\gamma)h^\dagger}}{4\gamma^2} \mathbb{E}_P [E_{t_n, n^\dagger h^\dagger} + h^\dagger \delta_\infty^2 + L_p^2 F_{t_n, n^\dagger h^\dagger}] \\
& \quad + \frac{5L_s^2 h^\dagger C_{K^\dagger}}{2\gamma} \left(h^{\dagger 2} e^{(3+\gamma)h^\dagger} (3\gamma d + M_s^2) + h^\dagger e^{2h^\dagger} d \right) \\
& \lesssim D_{\text{KL}}(\pi_{t_n, (n-1)h^\dagger} \| \tilde{\pi}_{t_n, (n-1)h^\dagger}) \\
& \quad + 5 \frac{2\gamma + L_s^2 C_{K^\dagger} 5h^{\dagger 2} e^{(3+\gamma)h^\dagger}}{4\gamma^2} \left(h^\dagger (L_s^2 + L_p^2) W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + h^\dagger \delta_\infty^2 + dh^\dagger \epsilon^{\dagger 2} \right) \\
& \quad + \frac{5L_s^2 h^\dagger C_{K^\dagger}}{2\gamma} \left(h^{\dagger 2} e^{(3+\gamma)h^\dagger} (3\gamma d + M_s^2) + h^\dagger e^{2h^\dagger} d \right) \\
& \lesssim D_{\text{KL}}(\pi_{t_n, (n-1)h^\dagger} \| \tilde{\pi}_{t_n, (n-1)h^\dagger}) + h^\dagger W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + h^\dagger \delta_\infty^2 + dh^\dagger \epsilon^{\dagger 2} + e^{-K^\dagger} h^{\dagger 2} d,
\end{aligned}$$

and then sum over n to obtain

$$\begin{aligned}
& D_{\text{KL}}(\pi_{t_n, T^\dagger} \| \tilde{\pi}_{t_n, T^\dagger}) = D_{\text{KL}}(\pi_{t_n, N^\dagger h^\dagger} \| \tilde{\pi}_{t_n, N^\dagger h^\dagger}) \\
& \lesssim D_{\text{KL}}(\pi_{t_n, 0} \| \tilde{\pi}_{t_n, 0}) + N^\dagger h^\dagger W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + N^\dagger h^\dagger \delta_\infty^2 + dN^\dagger h^\dagger \epsilon^{\dagger 2} + e^{-K^\dagger} N^\dagger h^{\dagger 2} d \\
& = T^\dagger W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + T^\dagger \delta_\infty^2 + dT^\dagger \epsilon^{\dagger 2} + e^{-K^\dagger} T^\dagger h^\dagger d.
\end{aligned}$$

Then, it is straightforward to see that when the following order of the parameters holds

$$\begin{aligned}
& T^\dagger = \mathcal{O}(1), \quad h^\dagger = \Theta(1), \quad N^\dagger = \mathcal{O}(1), \\
& \epsilon^\dagger = \Theta(d^{-1/2}\delta), \quad M^\dagger = \mathcal{O}(d^{1/2}\delta^{-1}), \quad K^\dagger = \mathcal{O}(\log(d\delta^{-2}))
\end{aligned}$$

and $\delta_\infty \leq \delta$, we have

$$D_{\text{KL}}(\pi_{t_n, T^\dagger} \| \tilde{\pi}_{t_n, T^\dagger}) \lesssim W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + \delta^2.$$

□

Lemma C.18. Suppose $T^\dagger \lesssim L_p^{-1/2}$, then we have

$$\text{TV}(\pi_{t_n, T^\dagger}, \tilde{p}_{t_{n+1}}) \leq \sqrt{D_{\text{KL}}(\pi_{t_n, T^\dagger} \| \tilde{p}_{t_{n+1}})} \lesssim \frac{1}{L_p^{\frac{1}{4}} (T^\dagger)^{\frac{3}{2}}} W_2(\pi_{t_n, 0}, \tilde{p}_{t_{n+1}}) \lesssim W_2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}).$$

Proof. A complete proof of the Lemma above is presented in [111, Lemma 9], which is derived based on [145, Corollary 4.7 (1)]. □

C.4 Overall Error Bound

We are now ready to prove Theorem 3.5.

Proof of Theorem 3.5. Notice that the interpolating corrector process $(\tilde{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}, \tilde{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger})$ is constructed to follow the same dynamics as the auxiliary corrector process $(\hat{\mathbf{u}}_{t_n, n^\dagger h^\dagger, \tau^\dagger}, \hat{\mathbf{v}}_{t_n, n^\dagger h^\dagger, \tau^\dagger})$ in the corrector step. Therefore, we have by data processing inequality that

$$\mathrm{TV}(\hat{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \tilde{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}) \leq \mathrm{TV}(\hat{\pi}_{t_n, 0}^{\tilde{\mathbf{u}}}, \tilde{\pi}_{t_n, 0}^{\tilde{\mathbf{u}}}) = \mathrm{TV}(\hat{q}_{t_n, h_n}, \tilde{q}_{t_n, h_n}), \quad (\text{C.32})$$

and again, since the interpolating predictor process $\tilde{\mathbf{y}}_{t_n, n^\dagger h^\dagger}$ is constructed to follow the same dynamics as the auxiliary predictor process $\hat{\mathbf{y}}_{t_n, n^\dagger h^\dagger}$ in the predictor step, we further have by data processing inequality that

$$\mathrm{TV}(\hat{q}_{t_n, h_n}, \tilde{q}_{t_n, h_n}) \leq \mathrm{TV}(\hat{q}_{t_n, 0}, \tilde{q}_{t_n, 0}) = \mathrm{TV}(\hat{q}_{t_n}, \tilde{p}_{t_n}). \quad (\text{C.33})$$

Furthermore, applying triangle inequality, Pinsker's inequality along with Theorem C.17 and Theorem C.7 proved above, we may upper bound the second term above as follows

$$\mathrm{TV}(\pi_{t_n, T^\dagger}, \tilde{\pi}_{t_n, T^\dagger})^2 \lesssim D_{\mathrm{KL}}(\pi_{t_n, T^\dagger} \| \tilde{\pi}_{t_n, T^\dagger}) \lesssim W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + \delta^2 \quad (\text{C.34})$$

Summarizing the above inequalities, we have

$$\begin{aligned} \mathrm{TV}(\tilde{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \tilde{p}_{t_{n+1}})^2 &= \mathrm{TV}(\tilde{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \pi_{t_n, T^\dagger}^{*, \mathbf{u}^*})^2 \\ &\leq \mathrm{TV}(\tilde{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \pi_{t_n, T^\dagger}^{\mathbf{u}})^2 + \mathrm{TV}(\pi_{t_n, T^\dagger}^{\mathbf{u}}, \pi_{t_n, T^\dagger}^{*, \mathbf{u}^*})^2 \\ &\leq \mathrm{TV}(\tilde{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \pi_{t_n, T^\dagger})^2 + \mathrm{TV}(\pi_{t_n, T^\dagger}, \pi_{t_n, T^\dagger}^*)^2 \\ &\leq \mathrm{TV}(\tilde{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \pi_{t_n, T^\dagger})^2 + \mathrm{TV}(\pi_{t_n, T^\dagger}, \tilde{p}_{t_{n+1}})^2 \\ &\lesssim W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) + \delta^2 + W_2^2(\tilde{q}_{t_n, h_n}, \tilde{p}_{t_{n+1}}) \\ &\lesssim de^{-K} + h_n^2 \delta_\infty^2 + d\epsilon^2 h_n^2 + \delta^2, \end{aligned} \quad (\text{C.35})$$

where the second last inequality is deduced from Theorem C.17 and Lemma C.18 and the last inequality is derived via Theorem C.7. Therefore, for any $n \in [0 : N - 1]$, applying triangle inequality along with data processing inequality (cf. Theorem A.1) yields

$$\begin{aligned} \mathrm{TV}(\hat{q}_{t_{n+1}}, \tilde{p}_{t_{n+1}}) &= \mathrm{TV}(\hat{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \tilde{p}_{t_{n+1}}) \\ &\leq \mathrm{TV}(\hat{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \tilde{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}) + \mathrm{TV}(\tilde{\pi}_{t_n, T^\dagger}^{\tilde{\mathbf{u}}}, \tilde{p}_{t_{n+1}}) \\ &\leq \mathrm{TV}(\hat{q}_{t_n}, \tilde{p}_{t_n}) + d^{1/2} e^{-K/2} + h_n \delta_\infty + d^{1/2} \epsilon h_n + \delta. \end{aligned} \quad (\text{C.36})$$

where the last inequality is derived by plugging in (C.32), (C.33) and (C.36). Applying Lemma A.9 and summing the inequalities above further give us that

$$\begin{aligned} \mathrm{TV}(\hat{q}_N, p_\eta) &= \mathrm{TV}(\hat{q}_N, \tilde{p}_N) \\ &\lesssim \mathrm{TV}(\hat{q}_0, \tilde{p}_0) + \sum_{n=0}^{N-1} \left(d^{1/2} e^{-\frac{K}{2}} + h_n \delta + d^{1/2} \epsilon h_n + \delta \right) \\ &\lesssim d^{1/2} e^{-T/2} + Nd^{1/2} e^{-K/2} + T\delta_\infty + d^{1/2} \epsilon T + \delta N. \end{aligned} \quad (\text{C.37})$$

By setting the parameters

$$\begin{aligned} T &= \mathcal{O}(\log(d\delta^{-2})), \quad h = \Theta(1), \quad N = \mathcal{O}(\log(d\delta^{-2})), \\ \epsilon &= \Theta\left(d^{-1/2} \delta \log^{-1}(d^{-1/2} \delta^{-1})\right), \quad M = \mathcal{O}(d^{1/2} \delta^{-1} \log(d^{1/2} \delta^{-1})), \quad K = \tilde{\mathcal{O}}(\log(d\delta^{-2})), \end{aligned}$$

and letting $\delta_\infty \lesssim \delta T^{-1} \lesssim \delta \log^{-1}(d\delta^{-2})$, we finally obtained the upper bound

$$\mathrm{TV}(\hat{q}_N, p_\eta)^2 \lesssim de^{-T} + N^2 de^{-K} + \delta^2 + d\epsilon^2 T^2 \leq \delta^2$$

as desired. \square

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We have carefully reviewed the abstract and introduction to ensure that they accurately reflect the paper's contributions and scope.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We have discussed the limitations of our work in Section 4 (Discussion and Conclusion).

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We have provided the full set of assumptions and complete proofs for all theoretical results in the paper.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: As a theoretical paper, we do not include experiments.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: As a theoretical paper, we do not include experiments.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: As a theoretical paper, we do not include experiments.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: As a theoretical paper, we do not include experiments.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: As a theoretical paper, we do not include experiments.

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We have reviewed the NeurIPS Code of Ethics and have ensured that our research conforms to it.

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: As a theoretical paper, we do not expect direct societal impacts of our work.

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: As a theoretical paper, we do not release data or models that have a high risk for misuse.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: As a theoretical paper, we do not use existing assets.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: As a theoretical paper, we do not introduce new assets.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: As a theoretical paper, we do not involve crowdsourcing nor research with human subjects.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: As a theoretical paper, we do not involve crowdsourcing nor research with human subjects.